# SIAM Conference on
# Parallel Processing for Scientific Computing
## (PP24)

**March 5–8, 2024**
**Lord Baltimore Hotel, Baltimore, Maryland, U.S.**

*This document was current as of February 26, 2024.*
*Abstracts appear as submitted.*

## IP1
### Frontier: The World's Most Powerful Computer for Science

The first exascale computer, called Frontier, was allocated for computational science at Oak Ridge National Laboratory in 2023. This unique scientific instrument is the culmination of more than a decade of concerted effort. I will relate a bit of the history of hybrid-node computing at the Oak Ridge Leadership Computing Facility (OLCF) and how Frontier represents the latest iteration of that approach. Some details of Frontier's architecture will be discussed, including an overview of the new AMD GPUs that provide the bulk of the computational power for Frontier. Finally, we will take a look at some problems that will benefit from the increased capability at exascale and I will convey some lessons learned from attacking these problems on this machine.

Bronson Messer
Oak Ridge National Laboratory
bronson@ornl.gov

## IP2
### The Power of Less: Harnessing Sparsity for Performance Optimization

Sparse matrix computations are fundamental to scientific computing and data analytics applications, such as computer graphics and machine learning. Sparsity leads to irregular memory accesses that pose challenges for code optimization. While specialized libraries can accelerate sparse computations, costly manual tuning is required for each application and architecture, reducing programmer productivity. Additionally, in machine learning, the unstructured sparsity patterns of deep learning models have rendered many of these libraries useless, prompting practitioners to use dense routine calls. Automation approaches such as compilers and runtime systems provide portability and ease of programming, but efficiently optimizing sparse codes remains a challenge due to access pattern irregularities. In this talk, I will introduce our work on building compilers and automation frameworks to accelerate sparse numerical kernels. I will present a class of inspection strategies that automatically analyze information such as matrix sparsity patterns and the numerical methods' properties, to generate optimized sparse codes. By leveraging formal verification and dynamic checks, we reduce reliance on domain-specific assumptions and enable correct and fast transformation of sparse codes. Additionally, I will discuss algorithmic modifications that make machine learning and graphics applications more amenable to sparse code specialization and the use of sparse compilers.

Maryam Mehri Dehnavi
University of Toronto
mmehride@cs.toronto.edu

## IP3
### Prospects for Efficient General-purpose Algebraic Solver Libraries for Multi-node GPU System

Multi-node GPU systems introduce difficulties for highly efficient general-purpose sparse algebraic solver software libraries. Open-source "subroutine libraries" (which include object-oriented libraries where the subroutines are class methods) provide efficient sparse algebraic solvers to multitudes of applications for the era of CPU-based high-performance computing. In such libraries, software implementations of complex algorithms are decomposed into a series of subroutine calls to a set of well-defined, relatively simple, separately-optimized computational kernels. The Basic Linear Algebra Subprograms (BLAS) provides one such decomposition for dense matrix computations. The decomposition of the library into a finite number of basic kernels limits the complexity of the software, making it possible to develop, maintain, and upgrade libraries that implement many complex algorithms. However, simple ports of subroutine libraries to GPUs tend to be inefficient due to large "kernel launch" times inherent to GPUs. Synchronously calling a series of GPU kernels from the CPU cannot efficiently utilize the GPU (even when all the data remains on the GPU). There are two techniques to limit the large launch time penalty: subroutine (kernel) fusion and asynchronous programming. Both come with the cost of additional software complexity, making managing general-purpose software libraries more difficult. The subroutine library paradigm for algebraic solvers may be reaching the end of its lifetime, but what will replace it?

Barry Smith
Simons Foundation and Flatiron Institute, U.S.
bsmith@petsc.dev

## IP4
### Towards Zero-waste Computing Through Co-design: The Case of Graph Processing

Graphs are universal abstractions, often used to represent concepts, objects, or individuals and the relations between them. When graphs are used to solve real-life problems—like the analysis of emerging inter-molecular interactions, the efficiency and feasibility of logistics networks, the reliability of support networks, or the safety of social networks—they quickly become massive in scale and/or complexity, and thus require massive computational resources to be processed. In turn, such large-scale processing of graph workloads raises efficiency and sustainability concerns: the irregular, data-intensive nature of graph processing algorithms often leads to significant computing waste. To alleviate these concerns, we propose a co-design methodology to enable the selection of efficient graph-processing algorithms *and* their effective deployment on suitable infrastructure. Our approach relies on design-space exploration, driven by efficient search methods and compositional performance models. We demonstrate our approach for solving a graph query that combines different types of graph processing or large-scale data, and show how our approach can significantly reduce the workload's energy consumption, while increasing resource utilization with limited performance impact. We conclude that our co-design approach is a leap forwards towards zero-waste computing for graph processing, and beyond.

Ana Lucia Varbanescu
University of Twente
a.l.varbanescu@utwente.nl

## IP5
### Performance Portability in the Age of Diverse Exascale Architectures

A diverse range of computer architectures have been explored to achieve Exascale-class supercomputers, including CPUs and GPUs from all the main vendors (AMD, Intel and NVIDIA), as well as the CPU-focused Fugaku system at RIKEN, which utilises Arm-based wide-vector CPUs from Fujitsu. For scientific applications to be able

to exploit as many of these systems as possible, they need to develop performance portable techniques. In this talk we will explore the latest achievements in performance portability. We will also discuss the latest developments in the path to Exascale in the UK.

Simon McIntosh-Smith
University of Bristol
cssnmis@bristol.ac.uk

## IP6
## Challenges of Scaling Deep Learning on HPC Systems

Machine learning algorithms, and training LLMs in specific, are becoming one of the main workloads running on HPC systems. More so, the scientific computing community is increasingly adopting modern deep learning approaches to their workflows. When HPC practitioners attempt to scale a typical HPC workload, they are mostly challenged by a particular bottleneck. Scaling deep learning, on the other hand, can be challenged by different bottlenecks: memory capacity, communication, I/O, compute etc. In this talk we give an overview of the bottlenecks in scaling deep learning, and highlight efforts in addressing some of those bottlenecks.

Mohamed Wahib
RIKEN Center for Computational Science
mohamed.attia@riken.jp

## IP7
## Do We Still Need Floating Point Arithmetic?

Block floating point (BFP) arithmetic is currently seeing a resurgence in interest for machine learning and AI applications. This is largely because it requires less power, less chip area, and is less complicated to implement in hardware compared to standard floating point arithmetic. This talk explores the application of BFP in scientific computing through mixed- and progressive-precision multigrid methods. Effectively, the use of BFP enables the solution of linear elliptic partial differential equations in energy- and hardware-efficient integer arithmetic. While most existing applications of BFP arithmetic tend to use small block sizes, the block size can be chosen to be maximal such that matrices and vectors share a single exponent for all entries. This is sometimes also referred to as a scaled fixed-point format. With this scheme some computations can be performed using as little as 4-bit integers while still reaching discretization-error-accuracy. For parallel implementations, normalization of intermediate results still poses a challenge, but for full multigrid it is possible to achieve this discretization-error-accuracy without any normalization steps. Ultimately, this leads to the broader question of whether we really need floating point arithmetic as much as we think we do?

Rasmus Tamstorf
Independent
rasmus.tamstorf@gmail.com

## SP1
## SIAM Activity Group on Supercomputing Best Paper Prize: Accelerating Sparse Iterative Solvers and Preconditioners Using RACE

The sparse matrix-vector multiplication (SpMV) kernel is a key performance-limiting component of numerous algorithms in computational science. Despite the kernels apparent simplicity, the sparse and potentially irregular data access patterns of SpMV and its intrinsically low computational intensity have been challenging the development of high-performance implementations of sparse algorithms over decades. In this talk, we present methods to increase the computational intensity and thereby accelerate the performance of SpMV kernels. The method is based on the concept of levels as developed in the context of our RACE library framework. We demonstrate that one can typically achieve a speedup of 1.5-4x on a single modern Intel or AMD multicore chip for symmetric SpMV and matrix power kernels using this level-based approach. After briefly introducing the optimization strategy, we apply these optimized kernels in iterative solvers. To this end, we discuss the coupling of the RACE library with the Trilinos framework and address the application to communication-avoiding s-step Krylov solvers, polynomial preconditioners, and algebraic multigrid (AMG) preconditioners. We then dive into the performance benefits and challenges of the RACE integration and show that our optimization produces numerically identical results and improves the total solver time by 1.3x - 2x.

Christie Alappat
Friedrich-Alexander-Universität Erlangen-Nürnberg,
Germany
christie.alappat@fau.de

Achim Basermann
German Aerospace Center (DLR)
Simulation and Software Technology
Achim.Basermann@dlr.de

Alan Bishop
Los Alamos National Laboratory
arb@lanl.gov

Holger Fehske
University of Greifswald, Germany
fehske@physik.uni-greifswald.de

Georg Hager
Friedrich-Alexander-Universität Erlangen-Nürnberg
georg.hager@fau.de

Olaf Schenk
Università della Svizzera italiana
olaf.schenk@usi.ch

Jonas Thies
German Aerospace Center (DLR)
j.thies@tudelft.nl

Gerhard Wellein
Friedrich-Alexander-Universität Erlangen-Nürnberg
gerhard.wellein@fau.de

## SP2
## SIAM Activity Group on Supercomputing Early Career Prize: Scalability and Productivity in Data-Intensive Biological Research on Massively Parallel Systems

The use of massively parallel systems continues to be critical for processing large volumes of data at an unprecedented speed and for scientific discoveries in simulation-based research areas. Today, these systems play a cru-

cial role in new and diverse areas of data science, such as computational biology, deep learning, and data analytics. Computational biology is a key area of the rapid growth of computing. The growing volume of data and increasing complexity have outpaced the processing capacity of single-node machines in these areas, making massively parallel systems an indispensable tool. The diverse and non-trivial challenges of parallelism in these areas require computing infrastructures that go beyond the demand of traditional simulation-based sciences. However, programming on high-performance computing (HPC) systems poses significant productivity and scalability challenges. It is important to introduce an abstraction layer that provides programming flexibility and productivity while ensuring high system performance. It is then important to map and plan the abstracted computation and communication to the underlying system to achieve optimal performance and guide the design of future large-scale architectures. As we enter the post-Moore's Law era, effective programming of specialized architectures is critical for improved performance in HPC. As large-scale systems become more heterogeneous, their efficient use for new, often irregular, and communication-intensive data analysis computation becomes increasingly complex. In this talk, we discuss how to achieve performance and scalability on extreme-scale systems while maintaining productivity for new data-intensive biological challenges, through an appropriate high-performance abstraction, namely the use of sparse matrices as well as the use of novel heterogeneous hardware.

Giulia Guidi
Cornell University
Lawrence Berkeley National Laboratory
gg434@cornell.edu

## SP3

**SIAM Activity Group on Supercomputing Career Prize: Tackling High Dimensional Problems Through Randomization and Communication Avoidance**

To Follow.

Laura Grigori
EPFL and PSI, Switzerland
laura.grigori@epfl.ch

## CP1

**A Mann-Type Iterative Scheme for Solving Variational Inclusion Problems with Applications**

Developing iterative schemes to tackle variational inclusion problems (VIP) has been a dynamic research field for several decades. In this presentation, we will provide an overview of the current algorithms designed to solve VIPs within a specific class. Additionally, we will introduce a novel and efficient algorithm based on the Mann method for VIP resolution. This approach incorporates inertial extrapolation terms and self-adaptive step sizes, offering improved convergence properties without the need for stringent conditions. This version is more flexible, featuring easily implementable criteria for the inertial factor and relaxation parameter. Finally, we will present some numerical experiments to demonstrate the practical implementation and comparative performance of the proposed algorithm.

Amara R. Eze
Morgan State University, Baltimore, MD

amara.eze@morgan.edu

Olaniyi Iyiola
Morgan State University
olaniyi.iyiola@morgan.edu

## CP1

**Utilizing Large Language Models for Disease Phenotyping in Obstructive Sleep Apnea**

Obstructive Sleep Apnea (OSA) affects millions in the U.S. and is linked to severe illnesses. However, despite its prevalence, there lacks understanding on how OSA interacts with common comorbidities and its additive risk of developing severe complications such as stroke and heart failure. This study focuses on utilizing Large Language Models (LLMs) to explore critical medical questions with the objective of obtaining actionable insights from clinical reports to model medical phenotypes and health metrics. We computed document-level embeddings for 331,793 discharge reports from the Medical Information Mart for Intensive Care IV (MIMIC-IV) database using the Perlmutter supercomputer. Twelve clinical LLMs were evaluated. Embeddings were clustered using K-Means. The purity of clusters was measured by Shannon entropy. Additionally, we assessed the quality of each model by visualizing its latent space with UMAP and manually reading sampled clinical text from clusters of interest. Clinical notes sampled from clusters of high rates of OSA patients with heart failure describe admissions of patients with a history of OSA, dyspnea as a chief complaint, and no prescription/adherence to Continuous Positive Airway Pressure (CPAP) treatment. This work serves to gain an understanding of how LLMs organize clinical information and if relations between latent representations capture phenotypes for feature discovery, which can be used to optimize the allocation of costly CPAP treatment.

Ifrah Khurram
San Juan Bautista School of Medicine
ifrahk@sanjuanbautista.edu

Rafael Zamora-Resendiz, Destinee S. Morrow
Hood College
Lawrence Berkeley National Laboratory
rzamoraresendiz@lbl.gov, dmorrow@lbl.gov

Silvia N. Crivelli
Lawrence Berkeley National Laboratory
University of California, Davis
sncrivelli@lbl.gov

## CP1

**Quantum-Powered Computational Multiphysics and Multiscale Modeling and Simulation**

The emergence of quantum computing technology has opened up new avenues in modeling and simulation in science and engineering. Traditional computational models and simulations have limitations in terms of accuracy, speed, and complexity. In this session, we discuss the potential of quantum computing in addressing these limitations, primarily from a computational science and engineering perspective, including computational fluid dynamics, finite element analysis, combustion modeling, fluid-structure interaction, and computational aero-acoustics. In this contribution, we highlight how quantum computing, combined with appropriate algorithms, can help us achieve unprece-

dented accuracy, speed, and complexity in simulations, eventually enabling significant progress in the design, testing, and development of complex engineered systems such as airplanes, spacecraft, cars, and machines. Our research and scientific computing software propose a paradigm shift in RD and engineering, bringing them more within the digital scope and realizing a significant speed-up in the design and development process of complex engineered systems and the simulation of multiphysics phenomena. Overall, this work demonstrates the transformative potential of quantum computing in computational modeling and engineering and multiphysics and multiscale simulation.

Mohamed Labadi
University of Chlef
Frontiers Labs
mohamed.labadi2@gmail.com

Samir Abdelmalek
University of Medea
Complex Systems Labs
abdelmalek.samir@univ-medea.dz

Ali Dali
National Center for Renewable Energy Development (CDER)
a.dali@cder.dz

Abdelkader Krimi
Ecole Polytechnique Montreal, Montreal-Canada
abdelkader.krimi@polymtl.ca

## CP1

### Simulations of Human-Scale Mars Lander Descent Trajectories on Leadership-Class Architectures

Future manned missions to Mars will require the safe delivery of substantially larger payloads to the planet surface as compared to prior robotic missions. To sufficiently decelerate such vehicles prior to touchdown, retropropulsion technologies are being considered. Such approaches introduce large uncertainties in vehicle performance and controllability during the reentry phase of flight. In this work, we present simulations performed on the GPU-based Summit and Frontier leadership-class computing systems located at Oak Ridge National Laboratory, featuring a real-time coupling with a flight dynamics package executing remotely at NASA Langley Research Center. Results include both perfect-gas and reacting-gas physical modeling strategies to simulate interactions between the LOX/CH4 rocket engines and the $CO_2$ Martian atmosphere.

Gabriel Nastac, Zachary Ernst, Kevin Jacobson, William Jones, Aaron Walden, Eric Nielsen
NASA Langley
gabriel.c.nastac@nasa.gov, zachary.ernst@nasa.gov, kevin.e.jacobson@nasa.gov, w.t.jones@nasa.gov, aaron.c.walden@nasa.gov, eric.j.nielsen@nasa.gov

Hayden Dean, Alexandra Hickey
Georgia Tech
hdean@gatech.edu, ahickey7@gatech.edu

Patrick Moran
NASA Ames

patrick.moran@nasa.gov

## CP2

### An AMG Reduction Framework for Domains with Symmetries

Divergence constraints are present in the governing equations of many physical phenomena, and they usually lead to a Poisson equation whose solution is one of the most challenging parts of scientific simulation codes. Algebraic Multigrid (AMG) is probably the most powerful preconditioner for Poisson's equation, and its effectiveness results from the complementary roles played by the smoother, responsible for damping high-frequency error components, and the coarse-grid correction, which in turn reduces low-frequency modes. This work presents a more compute-intensive variant of AMG. It arises from leveraging spatial symmetries, often present in academic and industrial configurations, to impose a consistent ordering giving rise to a multigrid reduction framework. In particular, we introduce an aggressive coarsening to the top level of the multigrid hierarchy, reducing the setup, memory footprint and application costs of the top-level smoother. Numerical experiments leveraging reflection and translational symmetries on CFD and structural mechanics problems will be presented at the conference.

del Alsalti-Baldellou
Polytechnic University of Catalonia
adel.alsalti@upc.edu

Carlo Janna
Dept. ICEA - University of Padova
carlo.janna@unipd.it

Andrea Franceschini, Gianluca Mazzucco
University of Padova
andrea.franceschini@unipd.it, gianluca.mazzucco@dicea.unipd.it

Xavier lvarez-Farré
SURF
xavier.alvarezfarre@surf.nl

F. Xavier Trias
Polytechnic University of Catalonia
francesc.xavier.trias@upc.edu

## CP2

### Domain Decomposition Methods for Neural PDE Solvers

Neural PDE solvers (NPS) offer notable advantages in approximating PDEs, including improved performance and stability with larger time steps in contrast to traditional numerical simulations. However, they face two key challenges that hinder their real-world application: the data-driven dilemma and geometric generalization. The data-driven dilemma stems from the need for extensive training data. Creating such datasets often relies on high-performance numerical solvers. Consequently, NPS tend to be evaluated on simplified problems, constrained to single physics, that are already solved efficiently with numerical methods. The geometric generalization dilemma arises from the need for PDE solvers to be used on complex and unique geometries–an area where NPS struggle. To address these challenges, we adapt domain decomposition methods (DDM) to NPS. Our approach enables mod-

els trained on small-scale problems to be applied to significantly larger problems without retraining. To emulate a Schwarz method with coarse correction, we incorporate graph neural networks (GNN) that propagate information globally. We also extend optimized Schwarz methods to NPS by learning transmission conditions during training. These techniques facilitate scalable information propagation across subdomains, ensuring fast approximation of the global solution. By combining DDM with GNNs, which excel in solving subdomains of arbitrary shape, we enable NPS to solve domains of varying size and complex geometry.

Arthur P. Feeney
University of California Irvine
afeeney@uci.edu

Aparna Chandramowlishwaran
UC Irvine
amowli@uci.edu

### CP2

**Trilinos Linear Solver Readiness for Modern Supercomputers**

The Trilinos project is a collection of high performance numerical libraries that underpin a wide range of applications. Among the offerings are a suite of linear solvers that run on a variety of computer architectures. In this talk, we present a "state of the union" assessment of Trilinos linear solvers with respect to current and emerging supercomputer platforms. In particular, we will describe the current readiness of Trilinos multigrid solvers and associated capabilities various GPU and CPU platforms. We will finish with an outline of ongoing and future plans.

Jonathan J. Hu
Sandia National Laboratories
Livermore, CA 94551
jhu@sandia.gov

### CP2

**Impact of Problem Structure on the Parallel Performance of Schur-Complement Decomposition Approaches for Large-Scale Nonlinear Optimization in Epidemic Inference**

The solution of large-scale nonlinear optimization problems using interior-point methods requires the solution of a sparse symmetric indefinite system at each iteration. For structured optimization problems, this linear system has block-bordered diagonal structure, which can be exploited for parallelization by employing a Schur-complement decomposition. The size and sparsity of the resulting Schur complement is problem specific, and in order to achieve good parallel speedup, different strategies must be adopted to process it. We apply both iterative and direct methods and investigate the parallel performance of the resulting Schur complement decomposition schemes for large-scale nonlinear parameter estimation problems in epidemic inference. Here, parameters of a national-scale, spatio-temporal model for epidemic diseases are fitted to county-level case data collected during the COVID-19 pandemic. We conduct our numerical experiments using parapint, a package to develop parallel algorithms for nonlinear optimization problems in Python. SNL is managed and operated by NTESS under DOE NNSA contract DE-NA0003525.

Laurens R. Lueg

Carnegie Mellon University
llueg@andrew.cmu.edu

Michael Bynum
Sandia National Labs
mlbynum@sandia.gov

Carl Laird, Lorenz Biegler
Carnegie Mellon University
claird@andrew.cmu.edu, biegler@cmu.edu

### CP2

**Gpu Performance of Algebraic Multigrid in Trilinos/MueLu**

We present results on the GPU performance of algebraic multigrid as implemented in Trilinos/MueLu with a focus on data migration and synchronization on GPU platforms. Our results will include both the setup and solve phases of the algorithm, on a single nodes as well as on multiple nodes.

Christopher Siefert
Sandia National Laboratories
csiefer@sandia.gov

### CP3

**Inversion of Time-Lapse Surface Gravity Data for Monitoring of 3D CO2 Plumes via Physics Informed Neural Networks**

We introduce two algorithms that invert simulated gravity data to 3D subsurface rock/flow properties. The first algorithm uses a supervised physics informed neural network, where we use a forward gravity model in our loss function. The second is an unsupervised physicsinformed approach that iteratively minimizes the difference between the gravity response of a predicted subsurface model and the original data. The target of these applications is the prediction of subsurface $CO_2$ plumes as a tool for monitoring $CO_2$ sequestration sites. Each proposed algorithm outperforms traditional inversion methods and existing datadriven deep learning approaches for the dataset at hand. Remarkably, our unsupervised model achieves better generalization to outof- distribution examples than its other deep learningbased counterparts. These results indicate that combining 4D surface gravity monitoring (low-cost acquisition) with physics informed deep learning techniques represents an effective and non-intrusive method for monitoring $CO_2$ storage sites.

Adrian Celaya
TotalEnergies EP Research and Technology, U.S. and Rice Univ
aecelaya@rice.edu

Bertrand Denel, Yen Sun, Mauricio Araya-Polo
TotalEnergies EP Research and Technology, U.S.
bertrand.denel@totalenergies.com,
yen.sun@totalenergies.com, mauricio.araya@shell.com

### CP3

**The Pele Simulation Suite for Reacting Flows at Exascale**

In this work, we present the Pele suite of software tools for compressible and incompressible reacting flows. The Pele suite leverages several different libraries, notably AMReX

and SUNDIALS, to achieve performance portability on heterogeneous computing architectures across the supercomputing landscape. The Pele suite is comprised of PeleC, a compressible reacting flow block-structured adaptive mesh refinement solver, PeleLMeX, a low-Mach number reacting flow block-structured adaptive mesh refinement solver, Pele- Physics, a library for transport, thermodynamics, finite rate chemistry, soot, spray and radiation physics. The objective of this paper is (i) to present the code development efforts necessary to achieve highly effective and scalable applications for exascale machines and (ii) to detail the performance results of the Combustion-Pele project applications on Oak Ridge National Laboratorys Frontier. We show good weak and strong scaling results for both PeleC and PeleLMeX up to more than 50 billion cells on more than 4096 Frontier graphics processing unit nodes. We also present a capability demonstration simulation of a dual-fuel pulse compression ignition engine (six adaptive mesh refinement levels, and 60 billion cells or 2.1 trillion degrees of freedom) on Frontier, to date one of the largest simulations performed on the first exascale-class supercomputer.

Marc T. Henry de Frahan
National Renewable Energy Laboratory, U.S.
Marc.HenrydeFrahan@nrel.gov

Lucas Esclapez
Lawrence Berkeley National Laboratory
lesclapez@lbl.gov

Jon Rood, Nicholas Wimer
National Renewable Energy Laboratory
jon.rood@nrel.gov, nicholas.wimer@nrel.gov

Paul Mullowney
Advanced Micro Device, U.S.
Paul.Mullowney@amd.com

Bruce A. Perry
National Renewable Energy Laboratory
bruce.perry@nrel.gov

Landon Owen
Sandia National Laboratories, U.S.
ldowen@sandia.gov

Hariswaran Sitaraman, Shashank Yellapantula, Malik Hassanaly
National Renewable Energy Laboratory
hariswaran.sitaraman@nrel.gov,
shashank.yellapantula@nrel.gov,
malik.hassanaly@gmail.com

Mohammad J. Rahimi
National Renewable Energy Laboratory, U.S.
mohammad.rahimi@nrel.gov

Michael Martin
National Renewable Energy Laboratory
michael.martin@nrel.gov

Olga A. Doronina
National Renewable Energy Laboratory, U.S.
olga.doronina@nrel.gov

Sreejith Appukuttan
National Renewable Energy Laboratory, Colorado
sreejith.nadakkalappukuttan@nrel.gov

Martin Rieth
Sandia National Laboratories
mrieth@sandia.gov

Wenjun Ge, Ramanan Sankaran
Oak Ridge National Laboratory, U.S.
gew1@ornl.gov, sankaranr@ornl.gov

Ann S. Almgren
Lawrence Berkeley National Laboratory
asalmgren@lbl.gov

Weiqun Zhang
Lawrence Berkeley National Laboratory
Center for Computational Sciences and Engineering
weiqunzhang@lbl.gov

John B. Bell
Lawrence Berkeley National Laboratory, U.S.
jbbell@lbl.gov

Ray W. Grout
National Renewable Energy Laboratory
ray.grout@nrel.gov

Marc Day
National Renewable Energy Laboratory, Colorado
marcus.day@nrel.gov

Jacqueline Chen
Sandia National Laboratories
jhchen@sandia.gov

**CP3**

**Scaling and Performance Portability of the Particle-in-Cell Scheme for Plasma Physics Applications Through Mini-Apps Targeting Exascale Architectures**

We perform a scaling and performance portability study of the electrostatic particle-in-cell scheme for plasma physics applications through a set of mini-apps we name Alpine, which can make use of exascale computing capabilities. The mini-apps are based on IPPL, a framework that is designed around performance portable and dimensionality independent particles and fields. We benchmark the simulations with varying parameters, such as grid resolutions (5123 to 20483) and number of simulation particles (109 to 1011), with the following mini-apps: weak and strong Landau damping, bump-on-tail and two-stream instabilities, and the dynamics of an electron bunch in a charge-neutral Penning trap. We show strong and weak scaling and analyze the performance of different components on several pre-exascale architectures, such as Piz-Daint, Cori, Summit, and Perlmutter. While the scaling and portability study helps to identify the performance critical components of the particle-in-cell scheme on the current state-of-the-art computing architectures, the mini-apps by themselves can be used to develop new algorithms and optimize their high performance implementations targeting exascale architectures.

Sriramkrishnan Muralikrishnan
Juelich Supercomputing Centre
Forschungszentrum Juelich GmbH
s.muralikrishnan@fz-juelich.de

Matthias Frey

University of St Andrews
mf248@st-andrews.ac.uk

Alessandro Vinciguerra
ETH Zurich
vincigua@student.ethz.ch

Micahel Ligotino
ETH Zurich, Switzerland
mic.ligotino@gmail.com

Antoine Cerfon
Courant Institute NYU
cerfon@cims.nyu.edu

Miroslav Stoyanov
Oak Ridge National Laboratory, U.S.
stoyanovmk@ornl.gov

Rahulkumar Gayatri
National Energy Research Scientific Computing Center,
U.S.
rgayatri@lbl.gov

Andreas Adelmann
Paul Scherrer Institut
andreas.adelmann@psi.ch

**CP4**

**Anderson Accelerated PMHSS for Complex-Symmetric Linear Systems**

This paper presents the design and development of an Anderson Accelerated Preconditioned Modified Hermitian and Skew-Hermitian Splitting (AA-PMHSS) method for solving complex-symmetric linear systems with application to electromagnetics problems, such as wave scattering and eddy currents. While it has been shown that the Anderson acceleration of real linear systems is essentially equivalent to GMRES, we show here that the formulation using Anderson acceleration leads to a more performant method. We show relatively good robustness compared to existing preconditioned GMRES methods and significantly better performance due to the faster evaluation of the preconditioner. In particular, AA-PMHSS can be applied to solve problems and equations arising from complex-valued systems, such as time-harmonic eddy current simulations discretized with the Finite Element Method. We also evaluate three test systems present in previous literature. We show that the method is competitive with two types of preconditioned GMRES, which share the significant advantage of having a convergence rate that is independent of the discretization size.

Måns I. Andersson
KTH Royal Institute of Technology
mansande@kth.se

Felix Liu
Raysearch Laboratories
KTH Royal Institute of Technology
felixliu@kth.se

Stefano Markidis
KTH Royal Institute of Technology, Sweden
markidis@kth.se

**CP4**

**Sequential and Shared-Memory Parallel Algorithms for Partitioned Local Depths**

In this work, we design, analyze, and optimize sequential and shared-memory parallel algorithms for partitioned local depths (PaLD). Given a set of data points and pairwise distances, PaLD is a method for identifying strength of pairwise relationships based on relative distances, enabling the identification of strong ties within dense and sparse communities even if their sizes and within-community absolute distances vary greatly. We design two algorithmic variants that perform community structure analysis through triplet comparisons of pairwise distances. We present theoretical analyses of computation and communication costs and prove that the sequential algorithms are communication optimal, up to constant factors. We introduce performance optimization strategies that yield sequential speedups of up to 29 over a baseline sequential implementation and parallel speedups of up to 26.2 over optimized sequential implementations using up to 32 threads on an Intel multicore CPU.

Aditya Devarakonda, Grey Ballard
Wake Forest University
devaraa@wfu.edu, ballard@wfu.edu

**CP4**

**Cholesky Factorization of Tile Low Rank Matrices on GPUs**

Tile low rank (TLR) representations of dense matrices partition them into blocks of roughly uniform size, where off-diagonal tiles are compressed and stored in low rank factorizations. They offer an attractive representation for many data-sparse dense operators that appear in practical applications, since substantial compression and a much smaller memory footprint can be achieved. Despite their utility, however, there are currently only a few high performance algorithms that can generate their Cholesky factorizations and operate on them efficiently, especially on GPUs. The difficulties in achieving high performance when factoring TLR matrices come from the expensive compression operations that must be performed during the factorization process and the irregular rank distribution of the tiles that requires an adaptive work pattern for the processing cores. In this work, we describe an algorithm that overcomes these limitations. Our algorithm has several new features. It always works in the compressed representation of the tiles. It compresses every tile in the output once only. It uses GEMM-rich adaptive randomized approximation for the compression. It also uses dynamic batched operations on the GPU to manage the irregular workload due to differing ranks among the output tiles. The resulting algorithm achieves substantial performance, as we demonstrate on sample matrices.

Wajih Boukaram
Lawrence Berkeley National Lab
wajih.boukaram@lbl.gov

Stefano Zampini
King Abdullah University of Science and Technology
stefano.zampini@kaust.edu.sa

George M. Turkiyyah
King Abdullah University of Science & Technology

(KAUST)
George.Turkiyyah@kaust.edu.sa

David E. Keyes
King Abdullah University of Science and Technology
(KAUST)
david.keyes@kaust.edu.sa

## CP5

### High-Throughput Scientific Computation with Heterogeneous Clusters: A Kitchen-Sink Approach Using the Actor Model

Scientific discovery has become increasingly reliant on high-throughput computation (HTC). HTC can be hindered, however, by issues such as a lack of accessibility to high-performance computing infrastructure or a lack of reliability (e.g., from volunteer computing). In this paper, we demonstrate how the actor model of concurrent computation offers the necessary tools to create customizable, robust, and scalable distributed HTC environments via a kitchen-sink approach, whereby all available computing resources are thrown at a given batch-based computation with the goal of maximizing throughput by maximizing accessibility. We assess the effectiveness of the kitchen-sink approach by applying it to a hydrological model, the Structure for Unification of Multiple Modeling Alternatives (SUMMA), to perform a simulation involving over half a million independent sub-simulations. We evaluate the proposed approach in two scenarios: one without node failures and one with multiple node failures. Our results affirm that the kitchen-sink approach not only successfully navigates these scenarios, but it also offers a novel and appealing approach to HTC.

Kyle Klenk
Department of Computer Science
University of Saskatchewan
kyle.klenk@usask.ca

Mohammad Mahdi Moayeri
University of Saskatchewan, Canada
mahdi.myr@gmail.com

Raymond J. Spiteri
University of Saskatchewan, Canada
Department of Computer Science
spiteri@cs.usask.ca

## CP5

### SYCL Compute Kernels for ExaHyPE

We discuss three SYCL realisations of a simple Finite Volume scheme over multiple Cartesian patches. The realisation flavours differ in the way how they map the compute steps onto loops and tasks: We compare an implementation that is exclusively using a sequence of for-loops to a version that uses nested parallelism, and finally benchmark these against a version modelling the calculations as task graph. Our work proposes realisation idioms to realise these flavours within SYCL. The results suggest that a mixture of classic task and data parallelism performs if we map this hybrid onto a solely data-parallel SYCL implementation, taking into account SYCL specifics and the problem size.

Tobias Weinzierl
Durham University
tobias.weinzierl@durham.ac.uk

Chung Min Loi
Durham University, United Kingdom
chung.m.loi@durham.ac.uk

Heinrich Bockhorst
Intel, Germany
heinrich.bockhorst@intel.com

## CP6

### Parallel Decomposition Strategies for a GPU Accelerated Continuum Kinetic Vlasov-Poisson Solver

Collisionless plasmas, in which mean free paths are longer than scales of interest, are important in astrophysical phenomena, electromagnetic propulsion, and fusion energy applications. Characterizing the behavior of these plasmas and their effect on transport properties requires kinetic theory, which treats each particle species as a probability distribution function in a six-dimensional position-velocity phase space. The evolution of high-dimensional distribution functions and associated low-dimensional electric fields is described by the Vlasov-Poisson system. Velocity moments connect low-dimensional and high-dimensional variables. In this talk we present strategies for effectively managing the coupling of the two spaces in a multi-GPU accelerated code. These include heuristics for partitioning the coupled structured domains, algorithms for performing a global reduction of the velocity moments and redistributing the results of the Poisson solve, and efficient GPU packing/unpacking algorithms for high-order finite-volume stencils. Strong and weak scaling of the algorithms will be compared, along with breakdowns of walltime spent on the various components in a timestep. Tests are performed using the GPU accelerated code VCK-GPU, a high-order accurate finite-volume Vlasov-Poisson solver. This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.

Andrew Ho, Genia Vogman
Lawrence Livermore National Laboratory
ho37@llnl.gov, vogman1@llnl.gov

## CP6

### Scaling Applications with Multi-level Parallelism via ML-guided Auto-Tuning

MPI+OpenMP is the de facto standard for scientific applications - or libraries for them, e.g., Kokkos, RAJA, HPX, Charm++ - run on supercomputers having nodes with heterogeneous processing units, e.g., CPU+GPU. While GPU parallelism libraries, e.g., CUDA, OpenMP-offload, are needed to harness computational power of GPUs, OpenMP is still key for harnessing multi-core parallelism on the CPU (host), especially with fatter nodes having heterogeneous sets of processing units, i.e., heterogeneous multi-xPUs. In this scenario, parameters of OpenMP must be tuned for node-level performance as well as across-node scalability with MPI. In this work, we present auto-tuning of OpenMP parameters *beneficial for scaling* MPI+OpenMP applications run on clusters having nodes with heterogeneous multi-xPUs. Our approach uses machine learning to guide tuning and pruning of a complex search space. We assess the benefits of our approach through experimentation on three benchmark codes relevant to scientific applications run with MPICH + LLVM's OpenMP + CUDA on DoE supercomputers. Our results suggest that our approach can improve application performance significantly

over the case of employing OpenMP auto-tuning without regard for MPI scalability.

Vivek Kale
Brookhaven National Laboratory
vivek.lkale@gmail.com

## CP6

### Scaling Hedgehog's Dataflow Graphs to Multi-Node GPU Architectures

Asynchuseronous task-based systems offer an array of abstractions to help us better utilize large-scale heterogeneous architecture. This paper extends the National Institute of Standards and Technology's Hedgehog [Alexandre Bardakoff et al., "Hedgehog: Understandable Scheduler-Free Heterogeneous Asynchronous Multithreaded Data-Flow Graphs," in 2020 IEEE/ACM 3rd Annual Parallel Applications Workshop: Alternatives To MPI+X (PAW-ATM), 2020, 1–15, (https://doi.org/10.1109/PAWATM51920.2020.00006.]

dataflow graph model to multinode GPU architectures [Nitish Shingde et al., "Extending Hedgehog's Dataflow Graphs toMulti-Node GPU Architectures," in Asynchronous Many-Task Systems and Applications, ed. Patrick Diehl et al. (Cham: Springer Nature Switzerland, 2023), 1–12.]. We present abstractions for doing workload balancing and data distributions to scale on larger clusters of nodes. The data decomposition proposed provides the flexibility to overcome memory constraints, whereas the data distribution adopted helps lower the overall internode communication volume to scale better. This newer approach is highlighted by using matrix multiplication as the driving vehicle. The performance results of this approach are measured against the leading libraries, SLATE and DPLASMA, for illustrative purposes only. This work aims to demonstrate that using general-purpose, high-level abstractions, such as Hedgehog's dataflow graphs, helps write scalable code without causing performance overhead.

Nitish Shingde
University of Utah
Scientific Computing and Imaging Institute
ngshingde@gmail.com

Martin Berzins
Scientific Computing and Imaging Institute
University of Utah
mb@sci.utah.edu

Timothy Blattner
National Institute of Standards and Technology
timothy.blattner@nist.gov

Walid Keyrouz
National Institute of Standards and Technology
Software and Systems Division
walid.keyrouz@nist.gov

Alexandre Bardakoff
National Institute of Standards and Technology
a.bardakoff@prometheuscomputing.com

## CP6

### Mixed Precision Arihtmetic to Accelerate a Hybrid Factorzation Solver for Sparse Matrices

A hybrid algorithm consists of decomposition of the sparse matrix into a union of moderate and hard parts during factorization with symmetric pivoting and generation of the Schur complement matrix for the hard part in higher precision by using the solution of the moderate part in lower precision. In precise, block GCR solver is used to generate the Schur complement matrix by solving moderate part in higher precision. We use LDU-factorized matrix in lower precision as a preconditioner for the iterative solver. In final, performing forward and backward substations for multiple RHS solution in higher precision with factorized matrices in lower precision is the essential part of the preconditioning procedure, where actual mixed precision arithmetic is necessary without type conversion of RHS data from higher to lower precision. Here, TRSM of BLAS level 3 is used for diagonal blocks and GEMM for off-diagonal ones for updating. These routines work with RHS in higher precision but coefficient matrix data in lower. The target problems are prepared in two different setting, to solve matrix with high condition number which exceeds the limit of double precision accuracy, where quadruple precision arithmetic is mandatory and one with moderate condition number which does not exceed the double precision accuacy, where faster computation with single precision is required. Double-double data is used as quadruple precision with acutal mixed precision for quadruple-double operations.

Atsushi Suzuki
Osaka University
atsushi.suzuki@cas.cmc.osaka-u.ac.jp

## CP6

### Gpu Accelerated Newton for Taylor Series Solutions of Polynomial Homotopies in Multiple Double Precision

A polynomial homotopy is a family of polynomial systems, typically in one parameter $t$. Our problem is to compute power series expansions of the coordinates of the solutions in the parameter $t$, accurately, using multiple double arithmetic. One application of this problem is the location of the nearest singular solution in a polynomial homotopy, via the theorem of Fabry. Power series serve as input to construct Pad approximations. Exploiting the massive parallelism of Graphics Processing Units (GPUs) capable of performing several trillions floating-point operations per second, the objective is to compensate for the cost overhead caused by arithmetic with power series in multiple double precision. The application of Newton's method for this problem requires the evaluation and differentiation of polynomials, followed by solving a blocked lower triangular linear system. Experimental results are obtained on NVIDIA GPUs, in particular the RTX 2080, RTX 4080, P100, and V100. Code generated by the CAMPARY software is used to obtain results in double double, quad double, and octo double precision. The programs in this study are self contained,

available in a public github repository under the GPL-v3.0 License.

Jan Verschelde
University of Illinois, Chicago, U.S.
janv@uic.edu

## CP7

### A Brick-based Hash-table Library for Structured-grid CFD on GPUs

This study is focused on porting a high-resolution computational fluid dynamics (CFD) code, featuring adaptive mesh refinement and mapped multi-block geometry, to the "HashBrick'library for execution on GPUs. Built on the concept of Bricks (https://bricks.run/index.html), which is a data layout and code generation framework enabling performance-portable stencil across a multitude of architectures, HashBrick introduces a new semi-structured organization of grid data optimized for sparsity and multilevel parallelism, providing improved data locality and cache reuse. The parallel performance will be demonstrated using a shockbox problem. The outcome of this study is establishing a performance model to characterize the optimal or sub-optimal parallel execution of the Hashbrick library in comparison to traditional structured-grid infrastructure. Future study will include more tests with the Hashbrick library on balancing the programmability, performance, portability.

Andrew Davis
University of Virginia
gtv6bk@virginia.edu

Stephen Guzik
Colorado State University
stephen.guzik@colostate.edu

Hans Johansen
Lawrence Berkeley National Laboratory
Computational Research Division
hjohansen@lbl.gov

Xinfeng Gao
University of Virginia
x.gao@virginia.edu

## CP7

### Coupling of Deep Learning and Bayesian Data Assimilation for High-Performance Computational Fluid Dynamics

Machine learning (ML) is now a core element in scientific computing (SC) and therefore potentially brings about new advances for computational fluid dynamics (CFD) modeling simulation of large-scale SC problems on high performance computing (HPC) architectures. The objective of this study is to investigate the performance in the coupling of ML and Bayesian data assimilation (DA) for CFD. DA is a method combining data (observations) with prior knowledge (e.g., CFD model) to improve the estimate of the distribution of the "true' state of a process (e.g., CFD simulation). Although ML and DA share some common mathematical and computational foundations, their difference and coupling with a large-scale SC code (e.g., CFD) merit an investigation of the performance and trade-offs. This study looks into the dataflow of each DA and ML process and assesses the performance of the coupling/feedback mechanism between DA, ML, and the "inner loop' CFD model.

Xinfeng Gao
University of Virginia
x.gao@virginia.edu

Youzuo Lin
Los Alamos National Laboratory
ylin@lanl.gov

## CP7

### Efficient Implicit Time-Stepping Schemes for the Incompressible Navier-Stokes Equations

We explore an algorithm for efficient solutions of three-dimensional incompressible flow problems as they appear in industrial applications like Formula 1 race cars. These complex geometries can be accurately represented using a spectral hp element method, however, the combination of high Reynolds number flows with high aspect ratio boundary layers strongly limits the stability of explicit schemes via a CFL condition. We are investigating an implicit time-stepping strategy to remove the time-step restriction and improve the robustness. The strategy is based on a linear-implicit Velocity-Correction scheme that decouples the monolithic system for pressure and velocity into one Poisson and three Advection-Diffusion-Reaction problems for each velocity component. Furthermore, the algorithm uses a linearised advection operator to remove the CFL restriction without coupling the system and, thus, stays linear in the time step which minimises any additional computational cost. We will present the algorithm itself and its performance on industrial problems with high Reynolds number turbulent flows.

Henrik Wüstenberg, Alexandra Liosi, Spencer Sherwin
Imperial College London
h.wustenberg@imperial.ac.uk, a.liosi22@imperial.ac.uk, s.sherwin@imperial.ac.uk

Joaquim Peiro
Dept of Aeronautics
Imperial College London, UK
j.peiro@imperial.ac.uk

David Moxey
Department of Engineering
King's College London
david.moxey@kcl.ac.uk

## CP7

### Performance Optimization of Automatic Differentiation for Unstructured-Grid CFD Applications on Emerging Heterogeneous Architectures

We share recent experiences performing kernel optimizations for typical motifs encountered in implicit methods

for Computational Fluid Dynamics (CFD) on unstructured grids. One particular focus area is Automatic Differentiation (AD) techniques aimed at obtaining Jacobians for complex physical models, often required for efficient solution methodologies, sensitivity analysis, error estimation, and uncertainty quantification. Such techniques present challenges for efficient implementation on Graphics Processing Units (GPUs), as complex physical models often involve many intermediate variables that create substantial register pressure and require significant memory traffic. We have developed an implementation that addresses these challenges and yields performance near the theoretical peak for todays state-of-the-art GPU architectures.

Eric Nielsen
NASA Langley
eric.j.nielsen@nasa.gov

Mohammad Zubair
Old Dominion University
zubair@cs.odu.edu

Desh Ranjan
Old Dominion University
Department of Computer Science
dranjan@cs.odu.edu

Aaron Walden, Gabriel Nastac
NASA Langley
aaron.c.walden@nasa.gov, gabriel.c.nastac@nasa.gov

Boris Diskin
National Institute of Aerospace
boris.diskin@nianet.org

## CP8
**Iteratively Decoupled Algorithms with a Parallel Stokes Subsolver for Biot's Model**

Biot's model is a multiphysics model which describes the interaction of a poroelastic material with its interstitial fluid flow. To find the numerical solution of Biot's model, we develop some iteratively decoupled algorithms that solve a reaction-diffusion subproblem for fluid pressure followed by solving a generalized Stokes subproblem. Our main algorithm is to parallel the generalized Stokes subsolver. Both analysis and experiments are given to demonstrate the effectiveness and efficiency of the algorithm. We theoretically show that the numerical solution of theiteratively decoupled algorithm converges to the solution of the coupled algorithm. Experimentally, we show that the parallel procedure does improve the efficiency and test the robustness of our algorithm with respect to physical and discretization parameters.

Mingchao Cai
Morgan State University
Morgan State University
Mingchao.Cai@morgan.edu

## CP8
**A Coupled Elliptic/Hyperbolic Adaptive Finite Volume Solver for High-Performance Heteroge-**neous Architectures

Coupled elliptic/hyperbolic systems of partial differential equations (PDEs) arise in a variety of science and engineering applications, including fluid solvers. Solving such systems on adaptive meshes allows for increased memory and compute efficiency but introduces additional complexities in implementation. We present our work on coupling the Hierarchical Poincar-Steklov (HPS) method, a fast and direct elliptic solver (Gillman and Martinsson, 2014), with the hyperbolic solvers implemented in the ForestClaw software (Calhoun and Burstedde, 2017). The underlying systems are solved on a dynamic quadtree mesh as implemented in the p4est library (Burstedde, et al., 2011). This includes targeting high performance, heterogeneous CPU/GPU architectures for distributed memory parallelism using our new workload sharing implementation of the HPS method. We will show scaling analysis on the Polaris petascale machine from Argonne National Laboratory and progress on our target application of a non-dispersive tsunami model.

Damyn M. Chipman, Donna Calhoun
Boise State University
DamynChipman@u.boisestate.edu,　　　donnacalhoun@boisestate.edu

## CP8
**Accelerating Time-Stepping Methods with Machine Learning Models and Hardware**

Many scientific simulations involve dynamical systems that are so complex and require such high fidelity that they become computationally impractical. Machine learning (ML) models offer cheaper and simpler ways to describe the dynamics of these systems at the cost of additional approximation errors. To overcome the individual limitations of the full and ML models, this talk presents new timestepping strategies based on multirate integrators that intelligently combine both models. The inexpensive ML model is integrated with a small timestep to guide the solution trajectory, and the full model is treated with a large timestep to occasionally correct for the ML model error and ensure convergence. Moreover, this approach can better utilize heterogenous computing resources by evaluating models on different hardware, e.g., the full model runs on a CPU while the ML model runs on a GPU or ML accelerator. Numerical experiment show that the hybrid integrators can be significantly more performant than using only the full or ML model for the integration.

Steven Roberts
Lawrence Livermore National Laboratory
roberts115@llnl.gov

## CP8
**An Efficient Parallelizable Exponential Time Differencing Method for Solving Reaction-Diffusion Equations**

There are different exponential time differencing (ETD) schemes for solving nonlinear differential equations such as ETD-Pade$^{'}$(1,1) or ETD-Crank-Nicolson scheme and ETD-Pade$^{'}$(0,2) scheme of second order convergent devel-

oped using well-known Pade′ approximation for the underlying matrix exponential. This study presents an efficient exponential time differencing scheme which utilizes the real distinct poles discretization for the matrix exponential called ETD-RDP. This scheme is proven to be second order convergent and most importantly enables parallel implementation of the resulting scheme. Comparison with second order ETD- Pade′ schemes (ETD-Crank-Nicolson scheme and ETD-Pade′(0,2) scheme) shows that the proposed ETD-RDP method is more efficient and accurate.

Madushi U. Wickramasinghe
Graduate Student (PhD)
Morgan State University
kewic1@morgan.edu

Olaniyi Iyiola
Morgan State University
olaniyi.iyiola@morgan.edu

**CP9**

**Ocean Models, But What If We Made Them 100× Faster?**

GPU computing has recently transformed high-performance computing (HPC) by providing substantial computational power and energy efficiency. However, ocean modeling, vital for climate and environmental research, remains primarily CPU-based, especially within the realm of unstructured models. Our work introduces a 2D/3D GPU ocean model based on the Discontinuous Galerkin (DG) method, overcoming challenges in adapting unstructured models to GPU architectures, traditionally optimized for structured memory access. Our contributions include memory efficiency optimizations and a novel "cell layout" memory structure. It strikes a balance between performance during the assembly, and speed of the resolution of the many banded linear systems that arise. Our GPU-accelerated model achieves an average 50× to 100× speedup over CPU counterparts, making a GPU equivalent to 1500-2500 CPU cores. These kind of speedups, while theoretically predicted, are rarely achieved in practice, especially for unstructured models and implicit processes. This research offers unprecedented computational efficiency for complex oceanic simulations, with direct ecological applications. It is already transforming our research group, where fine simulations that would previously require a year of computation can now be completed in just under a week.

Miguel De Le Court
UCLouvain
miguel.delecourt@uclouvain.be

Vincent Legat
Université catholique de Louvain
Belgium
vincent.legat@uclouvain.be

Jonathan Lambrechts
Université catholique de Louvain

jonathan.lambrechts@uclouvain.be

**CP9**

**Performance Characterization of MPAS-Ocean on Heterogeneous Supercomputers**

We present the performance of the Model for Prediction Across Scales (MPAS) Ocean model on Graphics Processing Unit (GPU)-based supercomputers. We evaluated the portability, functionality, and performance of the OpenACC port of MPAS-Ocean on multiple GPU architectures. Specifically, we have investigated performance on the United States Department of Energy's national laboratories' ecosystem supercomputers: Summit, Perlmutter-GPU, and Frontier, using Nvidia V100, A100, and AMD MI250X GPUs, respectively. We will present the computational performance of the ocean-sea ice coupled model as well as the standalone MPAS-Ocean. Moreover, the performance of the semi-implicit solver using linear iterative solvers in conjunction with linear algebra libraries for heterogeneous architectures, such as cuBLAS and hipBLAS, will be examined on different supercomputers. This activity is meant to understand deficiencies in the existing framework to inform the design of the next-generation ocean model titled OMEGA, the Ocean Model for E3SM Global Applications. We have already incorporated several lessons in the data structure and core framework design of OMEGA.

Youngsung Kim, Hyun Kang, Sarat Sreepathi
Oak Ridge National Laboratory
kimy@ornl.gov, kangh@ornl.gov, sarat@ornl.gov

**CP9**

**A New Coupler Infrastructure for E3SM Simulator**

Traditional coupled physics infrastructures for Climate models involve several steps in the simulation workflow to enable seamless transfer of information between components. The Model Coupling Toolkit (MCT) is one such software library that is used extensively in the E3SM and OASIS3 climate solvers. Since MCT does not store any detailed information about the underlying component meshes, the remapping weights to transfer data is computed using an offline process using tools such as ESMF, SCRIP or TempestRemap. In this talk, we present details about our ongoing work to integrate the MOAB mesh database within E3SM, with interfaces to TempestRemap to generate the remapping weights during the actual simulation, with zero offline requirements. Such a modified workflow maximizes the computational efficiency and increases scientific productivity, without sacrificing discretization accuracy of field data being transferred. The MOAB library has been interfaced with the HOMME (ATM), MPAS (OCN) and LND/RIVER models in E3SM to abstract out details of data transformation and migration between processes on large-scale machines. We will also present performance results of the new MOAB coupler, and provide comparative analysis against the existing E3SM-MCT coupler for problems of interest.

Vijay S. Mahadevan, Iulian Grindeanu, Robert Jacob
Argonne National Laboratory

mahadevan@anl.gov, iulian@anl.gov, jacob@anl.gov

## CP9

### Distributed and Gpu-Enabled Support Vector Machines for Wildfire Identification from Remote Sensing Data

We describe PermonSVM – a distributed memory parallel Support Vector Machine (SVM) implementation, built on top of the Portable, Extensible Toolkit for Scientific Computation (PETSc) – and its application to identification of wildfire-affected areas in large remote-sensing data sets. Automated machine learning approaches can enable more up-to-date and accurate tracking: the United States Monitoring Trends in Burn Severity ground truth layer we use is about two years behind and ignores small fires, which may be very important due to their frequency. Our aspirational goal is to enable synoptic mapping of fires across North America to enable better understanding of forest carbon balance. In our numerical experiments so far, we have achieved good classification performance and believe that our SVM classifier is effectively identifying smaller fires that are missing from the MTBS layers. Leveraging recent developments for GPU support in PETSc, we have achieved significant speedup for the SVM calculations on both NVIDIA and AMD GPU-powered supercomputers, and are actively working to further improve computational efficiency. Recent algorithmic improvements we are developing and evaluating include improved duality gap-based stopping criteria, novel approaches for stabilizing the underlying quadratic optimization algorithm, using Platt scaling to get probabilistic output from the classifiers, and GPU-friendly compact-dense formulations of quasi-Newton updates.

Richard T. Mills
Argonne National Laboratory
rtmills@anl.gov

Marek Pecha
Technical University of Ostrava
marek.pecha@vsb.cz

Zachary Langford, Jitendra Kumar
Oak Ridge National Laboratory
langfordzl@ornl.gov, jkumar@climatemodeling.org

David Horák
Technical University of Ostrava
david.horak@vsb.cz

## CP9

### Neural Network Prediction of Ocean Wave Behavior Using Frequency Domain Mapping

As ocean levels rise due to global warming, it is important to understand, visualize, and predict ocean wave behavior. Such predictions are an important component of emergency preparation. In recent years, key advances have been made in the application of neural networks toward weather related predictions. One neural network architecture, the Fourier Neural Operator (FNO), has shown great promise in image to image mapping where the output image repre-

sents the solution space of a physical system. In the current research, we use an FNO to map an ocean basin topography to the resulting wave height. We trained the FNO on a dataset generated from 1,600 simulations of the shallow water equations, where each simulation utilized a different arbitrarily-generated ocean basin topography. This type of simulation accurately depicts wave behavior away from the shoreline. Our setup is ideal for parallel computing. By parallelizing both the algorithm for the numerical method and the topography generation, we reduced computation time by 71.20%. Using parallel training of an ensemble of FNOs, we reduced the prediction MSE by 30.74%. Although we trained the FNO on a dataset of resolution 128 x 128, we validated the FNO on inputs of resolution 1024 x 1024 and received good results. For our future work, we plan to repeat this experiment using a more intensive approach such as particle-in-cell to expand the model dimensions from 2.5D to true 3D.

Pramya Surapaneni
High Technology High School
pramya.surapaneni@gmail.com

## CP10

### Advances in Data-Driven Solver Selection for Sparse Linear Matrices

Solving sparse linear systems is pivotal in diverse computational domains, and overall execution time of computations is heavily impacted by the efficient solution of such systems. Respective methods typically involve preconditioning an iterative solver; however, choosing optimal or sometimes even numerically stable combinations can be quite challenging. We discuss how to predict effective preconditioner and iterative solver combinations for any given sparse linear problem, through a combination of embedding and linear modelling techniques. We focus on determining useful system features and investigate different metrics to quantify the relative performance of the solvers across the SuiteSparse matrix collection on different architectures.

Hayden Liu Weng, Felix Dietrich
Technical University of Munich
h.liu@tum.de, felix.dietrich@tum.de

Hans Joachim Bungartz
Technical University Munich
bungartz@tum.de

## CP10

### Efficient Mesh Deformation Based on Randomized RBF Solvers

Mesh deformation methods have been widely used for the past decades in various fields such as fluid-structure interaction, aerodynamic shape optimization, unsteady and aeroelastic computational fluid dynamics. Among the existing methods, radial basis functions interpolation (RBF) is particularly suitable for unstructured mesh applications due to its simplicity and the high quality of the resulting mesh. Such approach requires solving dense linear systems, generally symmetric positive definite (SPD), which tends to be computationally expensive and memory demanding,

which is a major drawback when dealing with large-scale meshes. In this work, we aim to speed-up RBF mesh deformation procedure using methods coming from probabilistic linear algebra to solve the associated dense linear system. Indeed, such methods start to emerge into several fields, including numerical linear algebra and optimization, exploiting its spectral properties as an efficient alternative to reduce the complexity of solving large scale linear system. We will address the question of how to create an approximation to the RBF matrix by projecting the initial large scale operator onto a smaller subspace that exhibits specific properties such as a better sparsity and investigate various solving strategies. The proposed approach will be discussed on the basis of $2D$ and $3D$ applications.

Augustin Parret-Fréaud
Safran Tech, Rue des Jeunes Bois, Châteaufort
CS 80112 78772 Magny-Les-Hameaux
augustin.parret-freaud@safrangroup.com

Waël Bader
Safran Tech / CEMEF, Mines ParisTech, PSL, CNRS UMR 7635
wael.bader@safrangroup.com

Sebastien Da Veiga
ENSAI, CREST, CNRS UMR 9194
sebastien.da-veiga@ensai.fr

Youssef Mesri
MINES ParisTech, PSL University
youssef.mesri@mines-paristech.fr

## CP10

### Tensql: An SQL Database Built on Graphblas

Relational Database Management Systems (RDBMS) have been the most prominent form of database in the world for several decades. While relational databases are often applied within high-frequency/low-volume transactional applications such as website backends, the poor performance of relational databases on low-frequency/high-volume queries often precludes their application to big data analysis fields like graph analytics. This work explores the construction of an RDBMS solution that uses the GraphBLAS API to execute Structured Query Language (SQL) in an effort to improve performance on high-volume queries. Tables are redefined to be collections of sparse scalars, vectors, matrices, and more generally sparse tensors. The explicit values (nonzeros) in these sparse tensors define the rows and NULL values within the tables. A prototype database called TenSQL was constructed and evaluated against several SQL implementations including PostgreSQL. Preliminary results comparing the performance on queries common in graph analysis applications offer performance improvements as high as 1,400x over PostgreSQL for moderately sized datasets when returning results in a columnar format.

Jon P. Roose
Sandia National Laboratories
jproose@sandia.gov

## MS1

### From Iteration Counting to Applications with pySDC

Many parallel-in-time (PinT) algorithms replace the serial and, in this regard, direct way of time stepping by algorithms that iterate on multiple time steps in parallel. Efficient parallel implementation is a demanding exercise and may exceed the scope of more math-based PinT projects. Instead, PinT research typically resorts to "counting iterations" of the algorithm as a means of measuring its performance independently of its actual implementation. However, this hides many non-negligible sources of computational cost, such as communication. If the true parallel efficiency of candidate algorithms is under consideration, this needs to be obtained in a more coherent way. We present here one prototyping library which aims to cover both aspects of PinT research: method development and fair efficiency testing. pySDC is a Python code that, while not providing production-level performance, allows users to detect pitfalls in parallel algorithms before committing to optimized implementations or rigorous mathematical analysis. Parallel efficiency can be estimated by comparing to various time-serial algorithms implemented in the same framework. The modular structure allows users to easily apply a new algorithm to a wide range of problems and configurations without awareness of all the intricacies of the code. pySDC is publicly hosted on GitHub and well tested and documented in order to provide new users with a smooth and rapid start.

Thomas Baumann
Forschungszentrum Juelich
Technische Universitaet Hamburg Harburg
t.baumann@fz-juelich.de

Robert Speck
Juelich Supercomputing Centre
Forschungszentrum Juelich GmbH
r.speck@fz-juelich.de

Thibaut Lunet
Hamburg University of Technology
thibaut.lunet@tuhh.de

Lisa Wimmer
University of Wuppertal
wimmer@math.uni-wuppertal.de

Ikrom Akramov
Institute of Mathematics
Hamburg University of Technology
ikrom.akramov@tuhh.de

## MS1

### MaMiCo: High-Performance Parareal Molecular-Continuum Flow Simulation

Coupled molecular-continuum flow simulations extend the spatial domain of molecular dynamics (MD) using a computational fluid dynamics (CFD) solver. For a wide range

of scientific applications, these multiscale methods are an important tool to go beyond the limits of computational demand of pure MD. Often a massively parallel approach is necessary to yield acceptable time to solution. However, the spatial scalability of these methods with serial time stepping is limited, impeding researchers wanting to gain benefits from recent exascale high-performance computing systems. Hence, Parallel-in-Time (PinT) methods offer a promising way to enhance the scalability. In this talk, we present recent advances of the Parareal based PinT implementation in our open source molecular-continuum coupling tool, MaMiCo. MaMiCo is a C++ framework for 3D flows designed modularly, i.e. to offer coupling methodologies, including PinT, independent of the underlying MD and CFD codes. We employ an additional second continuum solver as hydrodynamic predictor to supervise the microscopic system. It reuses a variety of methods originally developed for coupling in space, such as noise filtering algorithms, mass flux and momentum transfer operators, to enable time parallelization. Our results demonstrate that the approach significantly improves the scalability of modular high-performance coupled flow simulations, opening the door to novel applications emerging on the horizon.

Piet Jarmatz
Helmut Schmidt University
jarmatz@hsu-hh.de

## MS1

### Multiscale Parareal Algorithms for Long-Time Mesoscopic Simulation of Complex Fluids

We present a multiscale parallel-in-time algorithm in which a continuum-based solver supervises a mesoscopic simulation in time-domain. Using an iterative prediction-correction algorithm, the parallel-in-time mesoscopic simulation supervised by its continuum-based counterpart can converge exponentially over iterations. The results show that the supervised mesoscopic simulations of both Newtonian fluids and non-Newtonian bloods converge to reference solutions after a few iterations. Physical quantities of interest including velocity, wall shear stress and flowrate are computed to compare against those of reference solutions, showing a less than 1% relative error on flowrate in the Newtonian flow and a less than 3% relative error in the non-Newtonian blood flow. The proposed method is then applied to a large-scale mesoscopic simulation of microvessel blood flow in a zebrafish hindbrain. The 3D geometry of the vasculature is constructed directly from the images of live zebrafish under a confocal microscope. The time-dependent blood flow from heartbeats in this realistic vascular network of zebrafish hindbrain is simulated using dissipative particle dynamics as the mesoscopic model, which is supervised by a one-dimensional continuum-based blood flow model in multiple temporal sub-domains. The computational analysis shows that the resulting microvessel blood flow converges to the reference solution after only two iterations.

Zhen Li
Clemson University

zli7@clemson.edu

## MS1

### Parallel-in-Time Software - Past, Present, Future

With million-way concurrency at hand, the efficient use of modern HPC systems has become one of the key challenges in computational science and engineering. For the numerical solution of time-dependent partial differential equations, time-parallel methods have shown to provide a promising way to extend prevailing scaling limits of numerical codes. However, while ideas, algorithms and proofs of concept are very visible in the literature, the corresponding codes and software packages are often hidden or only a marginal note in the various publications. The number of actually visible, time-parallel application codes is still rather small. Interestingly, the same is also true for stand-alone parallel-in-time libraries, which are mainly used for showcasing performance or for testing new ideas. In this talk, based on a recent survey conducted within the parallel-in-time community we will give an overview of what is currently available and describe some of the issues arising when implementing parallel-in-time integration methods. We show how authors see their codes, what software engineering techniques they apply and whether the codes themselves are perceived as scientific output.

Robert Speck
Juelich Supercomputing Centre
Forschungszentrum Juelich GmbH
r.speck@fz-juelich.de

Daniel Ruprecht
Institute of Mathematics
Hamburg University of Technology
ruprecht@tuhh.de

Jannik Finck
Juelich Supercomputing Centre
Forschungszentrum Juelich GmbH
j.finck@fz-juelich.de

## MS2

### Generative AI LLMs for High-performance Computing

Generative AI large language models (LLMs) have taken the world by storm since the release of GPT-3 in November 2022. There is a fundamental need to understand how these technologies are redefining how we do computing and high-performance computing (HPC). We evaluate the generation of HPC math kernels of the well-known GPT-3 OpenAI Codex, via Copilot, and Llama-2 foundational models and compare their accuracy on different languages (Fortran, C++, Python, Julia) and programming models (OpenMP, CUDA, HIP, Kokkos, OpenACC). We share our thoughts on the accuracy and an update on other aspects to add safeguards to the code generated by the foundational LLMs in question.

William F. Godoy
Oak Ridge National Laboratory

godoywf@ornl.gov

## MS2
### Openmp Offloading to Nvidia BlueField Dpus

Barcelona Supercomputing Center (BSC) has recently released OpenMP support for NVIDIA BlueField DPUs. This is an opportunity for the HPC community to leverage DPU features for a wider range of existing and emerging applications. In this talk I will introduce the current status of this support, its internal design and features, and the roadmap for upcoming features, along with the analysis of the first performance evaluations.

Antonio J. Peña, Sergio Iserte
Barcelona Supercomputing Center (BSC)
antonio.pena@bsc.es, sergio.iserte@bsc.es

## MS2
### Asynchronous Programming Models in Modern C++

Asynchronous many-task runtime systems have shown promising results lately in terms of helping to improve overall system utilization and parallel efficiency of scientific applications. One of the challenges in this domain is to design and maintain C++ APIs that a) are conforming to standard C++, b) unify local and remote operation (including working with accelerator devices), and c) are easy to use correctly. In this talk, we will highlight some of our recent results in this domain that demonstrate the viability of reaching all three of those goals.

Hartmut Kaiser
Louisiana State University
hkaiser@cct.lsu.edu

## MS3
### On Parallel Updates for Direct Factorization Methods

Computing has been disruptive to all scientific domains that increasingly rely on computational models and data to discover new knowledge and form decisions. With the explosion of Big Data, we are now faced with the ever-increasing size. The vast quantity, veracity, velocity, and variety of data are challenging classical high-performance numerical methods and software for extreme-scale computing. In this talk we address updating direct factorization methods, when being applied to sequences of large-scale problems where the systems are slightly varying from one step to the next step and recomputing the factorization would become too costly. Instead updating the existing factorization throughout the process of solving the systems is an attractive alternative. Depending on the problem, sometimes only parts of the factorization need to be updated. In other cases maybe global changes in the system require to update the whole factorization. These updates will be performed in parallel and, whenever possible, a large number of cores will be employed, even if eventually the factorization is not updated exactly but only by an iterative process. We will demonstrate the effectiveness of our approach for challenging large-scale application problems.

Olaf Schenk
Università della Svizzera italiana
olaf.schenk@usi.ch

Matthias Bollhoefer
TU Braunschweig
m.bollhoefer@tu-bs.de

## MS3
### Direct Sparse Solvers and Rank-Structured Preconditioners for Multi-GPU Systems

We present an approximate multifrontal sparse LU solver with block low rank compression. The solver runs on multi-GPU systems like Perlmutter and Frontier. The low rank compression is done with an adaptive randomized approximation, as implemented in the KBLAS library. The code also relies heavily on batched dense linear algebra routines from the MAGMA library. We discuss various pivoting and other numerical issues, resulting in a strong and robust preconditioner for a range of applications. We discuss applications like indefinite Maxwell, Navier-Stokes, and singularly perturbed reaction diffusion systems. We also compare with various other solvers and preconditioners.

Pieter Ghysels
Lawrence Berkeley National Laboratory
Computational Research Division
pghysels@lbl.gov

Lisa Claus
Lawrence Berkeley National Laboratory
lclaus@lbl.gov

## MS3
### Energy Earthshots, Heterogeneous Computing, and Linear Solvers

For the last two decades, large-scale irregular linear systems arising in engineering design computations have presented performance and scalability bottlenecks in multiple domain areas. These problems got renewed attention with the U.S. Department of Energy multi-billion-dollar investment in the Energy Earthshots initiative, which involves designing complex systems such as electrolyzers, gasifiers and flow batteries. While heterogeneous computing provides the computational power needed to support complex system design, effective hardware utilization requires rethinking the numerical linear algebra algorithms needed to solve the underlying sparse linear systems. Solutions currently used by domain experts rely heavily on sequential computations and, as such, severely restrict the amount of detail that can be included in a design computation. The lack of adequate numerical linear algebra capability was also a major obstacle for the Stochastic Grid Dynamics at Exascale (ExaSGD) subproject of the Exascale Computing Project (ECP). The research conducted within ECP, driven by power systems use cases from the ExaSGD subproject, has identified several promising directions on how to develop scalable direct (and direct-iterative) linear solvers that can handle very-sparse ill-conditioned systems. In this

talk we provide an overview of accomplishments achieved in this area within ECP and related activities at other programs.

Slaven Peles
Oak Ridge National Laboratory
peless@ornl.gov

Katarzyna Swirydowicz
Pacific Northwest National Laboratory
kasia.swirydowicz@pnnl.gov

Maksudul Alam
ORNL
alamm@ornl.gov

Koukpaizan Nicholson
Oak Ridge National Laboratory
koukpaizannk@ornl.gov

## MS3

### Portable Mixed-precision Algebraic Multigrid on GPUs

We present the GPU-native platform-portable algebraic multigrid (AMG) implementation in Ginkgo library that allows the user to use different precision formats from double, single, half, and bfloat16 precision for the different multigrid components. We uses an aggregation size 2 parallel graph match as the AMG coarsening strategy as demonstration in this presentation. The design provides a high level of flexibility in terms of configuring the bottom-level solver and the precision format for the distinct components of multigrid. We present the challenges, convergence, and performance results of mixed-precision AMG on GPUs.

Yu-Hsiang Tsai
Karlsruhe Institute of Technology, Germany
yu-hsiang.tsai@kit.edu

Natalie N. Beams
University of Tennessee
nbeams@icl.utk.edu

Hartwig Anzt
University of Tennessee, U.S.
hanzt@icl.utk.edu

## MS4

### Communication-Library Offload onto Smartnics (tentative)

This talk discusses design and implementation issues related to speeding up communication libraries using smart-NIC accelerators.

Richard Graham
NVIDIA
richgraham@nvidia.com

## MS4

### Emerging Use-Cases of SmartNICs for Scientific Computing

The emergence of SmartNICs has changed the field of networking by offloading complex processing and infrastructure tasks from the host CPU to the network interface. This talk focuses on the vast potential of SmartNICs, in particular in what the scientific community is interested on in terms of their use-cases. We will talk about how these smart network devices can improve various aspects of high performance computing, including communication offload, network security, storage and edge computing. The talk will also discuss research directions and challenges that need to be solved.

Oscar Hernandez
Oak Ridge National Lab
oscar@ornl.gov

## MS4

### Accelerating HPC and AI Applications by Offloading Computation and Communication to Smart-NICs

This talk presents new middleware (including an MPI implementation) for exploiting smartNIC devices, along with several preliminary application case studies.

Dhabaleswar K. Panda
The Ohio State University
panda@cse.ohio-state.edu

## MS4

### How Might Smart Networking Infrastructure Impact HPC Applications?

Smart networks, both in the form of smart NICs and smart switches, represent an interesting development in the hardware space and promise to significantly improve networking in the data center space, as well as support near storage compute. However, can they also improve classical HPC approaches, despite their limited compute capabilities? In the ScalNEXT project, funded by as part of the German ScalExa program for Exascale applications and systems, we are invistigating this questions, targeting both classic offload mechanisms as well as the use of smart networks to support program orchestration and coordination. In this talk, I will present this new project and its goal and will report on early developments and successes.

Martin Schulz
Technical University of Munich
schulzm@in.tum.de

Dennis-Florian Herr
echnical University of Munich
deflo.herr@tum.de

## MS5

### GPU-Accelerated Vecchia Approximations of Gaussian Processes using Batched Matrix Computations

Gaussian processes (GPs) are often used in geospatial

analysis but struggle with large data sets due to computational complexity. Specifically, calculating the log-likelihood function for geospatial data involves inverting a large dense covariance matrix. Researchers have been focusing on approximation methods for better scalability and accuracy, such as the Vecchia approximation. This study uses batched matrix computations on modern GPUs, leveraging the KBLAS linear algebra library to introduce a parallel Vecchia approximation. This approach significantly speeds up the log-likelihood function evaluation, achieving up to 700X, 1180X, and 900X faster performance on V100, A100, and H100 GPUs. It can handle up to 1 million locations on 80GB A100 and H100 GPUs with high accuracy. The study also tests the accuracy of this algorithm on real-world geospatial datasets, including soil moisture in the Mississippi Basin and wind speed in the Middle East.

Sameh Abdulah
Extreme Computing Research Center (ECRC)
King Abdullah University of Science and Technology
sameh.abdulah@kaust.edu.sa

**MS5**

**Multigrid Reduction Preconditioning for Coupled Subsurface Processes**

The simulation of subsurface fluid flow involves solving complex multi-physics problems with tightly coupled multiphase flow and transport. To address this dynamic interplay, fully implicit methods, known as monolithic approaches, are typically favored. However, these methods entail solving large, non-symmetric, indefinite, and highly ill-conditioned linear systems arising from discretization and linearization of governing balance equations, making them challenging. These problems become even more challenging when coupling fluid flow with mechanical processes like rock fracture and deformation. Hence, designing preconditioners that manage these physics couplings is essential for achieving rapid convergence. The advent of new high-performance computing (HPC) hardware featuring accelerators such as GPUs offers opportunities to enhance solution performance when effectively leveraged. This work presents our endeavors to develop scalable linear solver strategies tailored for subsurface flow applications, prioritizing efficiency on modern HPC systems. We introduce a framework centered on multigrid reduction (MGR), specifically designed for tightly coupled systems of partial differential equations (PDEs). We showcase its adaptability across various scenarios and demonstrate its efficiency and scalability when addressing substantial problems on contemporary HPC architectures.

Nicola Castelletto
Lawrence Livermore National Laboratory
Livermore, CA, USA
castelletto1@llnl.gov

Matteo Cusini
Lawrence Livermore National Laboratory
cusini1@llnl.gov

François Hamon
TotalEnergies E&P Research and Technology
francois.hamon@totalenergies.com

Daniel Osei-Kuffuor, Victor A. Paludetto Magri,
Randolph R. Settgast, William Tobin, Joshua A. White
Lawrence Livermore National Laboratory
oseikuffuor1@llnl.gov, paludettomag1@llnl.gov,
settgast1@llnl.gov, tobin6@llnl.gov, white230@llnl.gov

**MS5**

**Sparse Matrix As Computational Abstraction for Large-Scale Science**

In the last decade, advances in genome sequencing have flooded us with genomic data, leading to new computational approaches in bioinformatics. As genomics becomes increasingly important in the health and environmental fields due to the decreased cost of sequencing, the flood of data requires parallel processing, often on high-performance computing (HPC) systems. However, the complicated, irregular nature of genomic computation makes parallelism in distributed memory difficult. This is as true for other sciences as it is for data science. Our work demonstrates that it is possible to create highly parallel code for complex genomic computation by using an appropriate abstraction such as sparse matrices. In addition, we explore the use of specialized hardware such as GPUs. The result is high performance, high productivity, and an efficient computational pipeline.

Giulia Guidi
Cornell University
Lawrence Berkeley National Laboratory
gg434@cornell.edu

**MS5**

**Distributed Massive-Scale Sparse Inverse Covariance Estimation**

Across various scientific disciplines, sparse inverse covariance estimation is a popular tool for capturing the underlying dependency relationships in multivariate data. Unfortunately, most estimators are not scalable enough to handle the sizes of modern high-dimensional data sets (often terabytes). In statistics, a sparse inverse covariance matrix, or precision matrix, is a fundamental quantity for characterizing conditional dependencies in multivariate data. Estimating the precision matrix is complicated by various statistical and computational challenges. We describe a massive-scale sparse precision matrix estimation approach that can analyze data with 1 million variables, and application examples are discussed.

Sang-Yun Oh
University of California, Santa Barbara
syoh@ucsb.edu

**MS6**

**Lessons Learned and Portability Strategies for a Large-Scale Materials Science Codes to GPU Architecture**

Computational materials science codes have been and still are among the applications which mostly benefit from leadership class HPC facilities, motivating forefront software developments to keep up with the constant increase of com-

putational power offered by newer architectures. For over two decades, HPC software developers relied on advancements in massive parallelization to provide performance boosts. More recently, vendor-specific architectures (such as graphic processing units GPU) have emerged as the dominant HPC paradigm, forcing developers to actively maintain and optimize core compute kernels in a newer and more complex fashion. In this talk we focus on the lesson learned and the devised portability strategies on GPU architectures for a large-scale materials science code, namely BerkeleyGW. BerkeleyGW is a massively parallel software package employed to study the excited state properties of electrons in materials by using GW, Bethe-Salpeter Equation (BSE) methods and beyond. We discuss our software design to achieve true performance portability across various GPU vendor architectures by analyzing the performance of vendor specific programming models (CUDA and HIP) and the directives based portable counterparts (OpenACC and OpenMP-target). We highlight the challenges we encountered as well as the practices we found useful in the porting pipeline for the three major GPU vendor solutions, namely NVIDIA, AMD and Intel.

Mauro Del Ben
Lawrence Berkeley National Laboratory
mdelben@lbl.gov

### MS6

### Replacing Diagrammatic Quantum Monte Carlo with Tensor Trains

The simulation of strongly correlated quantum models is a significant challenge in modern condensed matter physics. Among the most successful methods are Monte Carlo techniques that accurately and reliably sample perturbative expansions to any order. However, the cost of obtaining high precision through these methods is high. Recently, tensor train decomposition techniques have been developed as an alternative to Monte Carlo integration. In this study, we apply these techniques to the single-impurity Anderson model at equilibrium by calculating the systematic expansion in power of the hybridization of the impurity with the bath. We demonstrate the performance of the method in a paradigmatic application, examining the first-order phase transition on the infinite dimensional Bethe lattice, which can be mapped to an impurity model through dynamical mean field theory. Our results indicate that using tensor train decomposition schemes allows the calculation of finite-temperature Green's functions and thermodynamic observables with unprecedented accuracy. The methodology holds promise for future applications to frustrated multi-orbital systems, using a combination of partially summed series with other techniques pioneered in diagrammatic and continuous-time quantum Monte Carlo.

Emanuel Gull, Andre Erpenbeck, Wei-Ting Lin, Thomas Blommel, Lei Zhang, Sergei Iskakov
University of Michigan
egull@umich.edu, aerp@umich.edu, wtlin@umich.edu, tblommel@umich.edu, lzphys@umich.edu, siskakov@umich.edu

Liam Bernheimer

Tel Aviv University
liamb@mail.tau.ac.il

Yuriel Nunez-Fernandez
CEA Grenoble
yurielnf@gmail.com

Guy Cohen
Tel Aviv University
guy.cohen@gmail.com

Olivier Parcollet
CCQ, Flatiron Institute
olivier.parcollet@cea.fr

Xavier Waintal
CEA Saclay
xavier.waintal@cea.fr

### MS6

### Parallel Coordinate Descent Full Configuration Interaction

We develop a mutli-threaded parallel coordinate descent full configuration interaction algorithm, for the electronic structure ground-state calculation in the configuration interaction framework. The algorithm solves an unconstrained nonconvex optimization problem, via a modified block coordinate descent method with a deterministic compression strategy. CDFCI captures and updates appreciative determinants with different frequencies proportional to their importance. We demonstrate the efficiency of the algorithm on practical systems.

Yuejia Zhang
Fudan University
yuejiazhang21@m.fudan.edu.cn

Weiguo Gao
Fudan University, China
wggao@fudan.edu.cn

Yingzhou Li
Fudan University
yingzhouli@fudan.edu.cn

### MS6

### Solving the Kadanoff-Baym Equation on GPUs

Computing the non-equilibrium Green's function of a quantum many-body system by solving the Kadanoff-Baym equations, a set of nonlinear integral-differential equations is challenging. We describe efforts to implement a KBE solver on multiple GPUs using both OpenACC and CUDA FORTRAN. To achieve high performance, several techniques (code restructuring, loop fusion and reordering, batching) are needed be achieve good performance for self-energy and collision integral calculations. We will give a few examples to demonstrate the effectiveness of these techniques.

Chao Yang
Lawrence Berkeley National Lab

cyang@lbl.gov

Jia Yin
Lawrence Berkeley National Laboratory
jiayin@lbl.gov

Khaled Z. Ibrahim
Lawrence Berkeley National Laboratory
Berkeley, CA, USA
kzibrahim@lbl.gov

Mauro Del Ben
Lawrence Berkeley National Laboratory
mdelben@lbl.gov

## MS7

**The Sharp Interface Simulation of Ternary Alloy Solidification with Melt Convection on Quadtrees**

We present a numerical method for the simulation of solidification of ternary alloy coupled with melt convection effects in 2D. A combination of spatially adaptive quadtree grids, level-set method, and sharp-interface numerical methods for imposing boundary conditions is used to accurately and efficiently resolve the complex behavior of the solidification front. The governing equations for ternary alloys lead to a non-linear system of coupled PDEs, and so we cannot use simple fixed point iterations for binary alloys. Instead, a modified Newton-type approach on an adaptive quadtree mesh is needed to capture the complex behavior of the solidification front, which is also coupled with the incompressible Navier-Stokes equations to include the fluid effects. We first verify the method with a numerical convergence test using a synthetic solution, to demonstrate that the pressure-free projection method accurately represents convection in the melt pool. Finally, we perform numerical experiments for directional solidification of a ternary alloy in a shear flow, and analyze the solutal segregation dependence on that and other processing conditions.

Rochi Chowdhury
University of California, Santa Barbara, U.S.
rochishnu00@ucsb.edu

Elyce Bayat
University of California, Santa Barbara
ebayat@ucsb.edu

Frederic Gibou
UCSB
fgibou@ucsb.edu

## MS7

**Is Higher Accuracy Free for Simulations with Interfaces?**

Many arguments for using GPUs for PDE simulations assume there will be benefits from high memory bandwidth and massive floating point capabilities. However, the fact is that the most-used algorithms have very low arithmetic intensity (AI, ratio of FLOPs/bytes), poor communication patterns, and dont take advantage of GPU compute capabilities. This is especially true for dynamic, multi-physics interfaces, since codes often use both adaptive grid refinement and local time stepping to try to improve poor space-time accuracy. I will present a few results suggesting a different goal: What accuracy can be achieved with the fewest resources and least coding effort? Testing this work-accuracy-effort tradeoff suggests that typical low-AI methods fail to achieve both efficient parallel scaling and high accuracy. In contrast, higher-order methods have lots of FLOPs, achieve better accuracy, faster, and lend themselves to abstractions that are easier to code. I will argue that better accuracy and scaling with GPUs *is* compatible with ease of coding, if we are willing to adopt higher-order algorithms.

Hans Johansen
Lawrence Berkeley National Laboratory
Computational Research Division
hjohansen@lbl.gov

## MS7

**Accurately Capturing Breakup Dynamics in Multiphase Flows Using Scalable Adaptive Algorithms**

We address a long-standing issue of artificial (numerical) breakup of fluid structures in the modeling of interface-resolved simulations of multiphase flows. The root of this artificial breakup is the comparable length scale between the mesh and the interfacial features in the flow, e.g., filaments, sheets, and droplets. Thus, the accuracy of these simulations is strongly contingent upon the finest mesh resolution used to represent the fluid-fluid interface structures. However, the increased resolution comes at a higher computational cost, even when using adaptive refinement strategies. We propose algorithmic advances that leverage adaptive octree-based meshing and aim to reduce the computational cost without compromising on the physics by selectively detecting key regions of interest (droplets/filaments) that require significantly higher resolution. We demonstrate the scaling of the framework up to 114,688 processes on Frontera. We deploy the framework to capture complex break-up dynamics of pulsed primary jet atomization with one of the highest resolved simulations enabled using this dynamic detection and targeted refinement approach. Our simulations reveal that pulsed jet atomization exhibits a complex cascade of break-up mechanisms involving sheet rupture and filament formation. The proposed approach opens up the ability to affordably perform resolved simulations for various multiphase flow phenomena exhibiting thin features and their breakup.

Makrand A. Khanwale
Stanford University, U.S.
khanwale@stanford.edu

Kumar Saurabh
Iowa State University
maksbh@iastate.edu

Masado Ishii
University of Utah
masado@cs.utah.edu

Hari Sundar

School of Computing
University of Utah
hari@cs.utah.edu

Baskar Ganapathysubramanian
Iowa State University
baskarg@iastate.edu

**MS8**

**Arachne: A Productive Massive-Scale Graph Analytics Framework**

Analyzing massive-scale networks, or graphs, presents significant challenges due to their vast size and the growing need for user-friendly and productive analytic tools. Such graphs hail from diverse fields such as cybersecurity and health representing anything from malicious network flows to neural connections. Users require simple and efficient tools to decipher hidden structures and derive insightful metrics from these massive graphs. Tools such as NumPy, Pandas, and NetworkX, allow users to analyze data in Python of limited size. Recently, Arkouda was released as a replacement at scale for array operations typically found in NumPy and Pandas. Arkouda allows for efficient execution of array tasks for datasets terabytes in size. Its operations on tabular data, like reductions and scans, differ from the often irregular nature of operations on graphs. Arachne, an extension to Arkouda, bridges this gap by seamlessly transforming Arkouda dataframes into graphs. Notably, Arachne can manage complex graph types, like property graphs, processing queries on graphs composed of billions of edges in seconds. Both Arkouda and Arachne harness the power of Chapel, a language crafted to elevate parallel programming by removing the need to manage inter-node communication in a computing cluster, with a Python API. This talk will provide a comprehensive insight into the capabilities of Arachne and Chapel's pivotal role in its design.

Oliver Alvarado Rodriguez
New Jersey Institute of Technology
oaa9@njit.edu

Zhihui Du
NJIT
zhihui.du@njit.edu

David A. Bader
New Jersey Institute of Technology
bader@njit.edu

**MS8**

**Towards Interactive HPC-Scale Property and Knowledge Graph Analytics with MetallData**

Scaling graph algorithms to ever-increasing data volumes has received tremendous attention from the HPC research community, primarily focusing on parallel, distributed, and/or accelerator-based graph algorithms. Venues such as the Graph500, GraphChallenge, and DIMACS challenges have focused research towards developing new processing capabilities. With advances in raw computing capabilities comes to need to enable exploratory interactive comput-

ing (e.g., Jupyter notebook driven) for HPC-scale analytics. In this talk, we present active research efforts investigating persistent data structures, and how they can enable HPC-scale interactive computing for graphs and other data-intensive applications. We will demonstrate how to drive an HPC graph analytics workflow from an interactive python environment that is extendable to other analytic domain areas. Additionally, present research investigating algorithms and techniques to scale graph analytics to some of the largest open-source datasets in distributed HPC environments, including the new AWS HPC-Cloud environment. We show that AWSs HPC environment is competitive with traditional HPC and is capable of operating on the full 14.3 billion comment Reddit graph using a suite of graph-based benchmarks.

Roger Pearce
Lawrence Livermore National Laboratory
rpearce@llnl.gov

Trevor Steil
LLNL
steil@gmail.com

Keita Iwabuchi
Lawrence Livermore National Laboratory, USA
iwabuchi1@llnl.gov

Peter Pirkelbauer
Lawrence Livermore National Laboratory
pirkelbauer2@llnl.gov

**MS8**

**Communication-Avoiding Algorithms for Full-Batch and Mini-Batch GNN Training**

Graph Neural Networks (GNNs) are powerful and flexible neural networks that use the naturally sparse connectivity information of graph data. GNNs represent this connectivity as sparse matrices, which have lower arithmetic intensity and thus higher communication costs compared to dense matrices, making GNNs harder to scale to high concurrencies than Transformers or fully-connected neural networks. We show that communication-avoiding matrix multiplication algorithms can accelerate GNN training compared to existing training methods for both minibatch and full-batch training. First, we introduce the first distributed GNN training system that supports multiple minibatch sampling algorithms, enabling GNN scales that were previously infeasible. Our approach is to express each step in GNN training with sparse matrix operations, and leverage communication-avoiding SpGEMM algorithms to scale training to hundreds of GPUs. Second, we show how full-batch training can scale out to hundreds of GPUs by leveraging communication-avoiding SpMM algorithms. We experiment on graph datasets with billions of edges on over a hundred GPUs to show the performance benefits of these algorithms.

Alok Tripathy, Katherine Yelick
University of California, Berkeley
Lawrence Berkeley National Laboratory
alokt@berkeley.edu, yelick@berkeley.edu

Aydin Buluc
Lawrence Berkeley National Laboratory
abuluc@lbl.gov

## MS9

### WarpX Experiences at Exascale

In July of 2022, the particle-in-cell code WarpX ran for the first time a laser-plasma simulation on nearly the full size of the first reported Exascale Supercomputer, Frontier (OLCF). Frontier's novel AMD GPU architecture complemented successful runs on the prior TOP1 supercomputers (Fugaku, A64FX CPUs) and Summit (Nvidia V100 GPUs), along with the latest NERSC Supercomputer (Perlmutter, Nvidia A100 GPUs). Looking back to 2016, the road to sustainably develop an application that has many numerical innovations was not yet clear. Would performance portability layers become general enough and provide low enough overhead to enable single-source development? At which point in the software stack should hardware-specialization routinely start? How to ensure low-entry burdens, extensibility and sustain numerical innovation with a large community? This posed a major challenge for design and maintainability. In this talk, we will show how integration with and contribution to the co-design framework AMReX, cross-expertise teamwork, effective challenging and leveraging software library collaborations and focused performance tracing over the project runtime enabled the WarpX application to achieve its science goals, and attract a large open source community sustaining it beyond the DOE ECP Project.

Axel Huebl
Lawrence Berkeley National Laboratory (LBNL), USA
axelhuebl@lbl.gov

Rémi Lehe
Lawrence Berkeley National Laboratory
ATAP
rlehe@lbl.gov

Olga Shapoval
LBNL
OShapoval@lbl.gov

Edoardo Zoni, Arianna Formenti, <u>Marco Garten</u>
Lawrence Berkeley National Laboratory
ezoni@lbl.gov, ariannaformenti@lbl.gov, mgarten@lbl.gov

Andrew Myers
Lawrence Berkeley National Lab
atmyers@lbl.gov

Weiqun Zhang
Lawrence Berkeley National Laboratory
Center for Computational Sciences and Engineering
weiqunzhang@lbl.gov

John B. Bell
Lawrence Berkeley National Laboratory, U.S.
jbbell@lbl.gov

Ann S. Almgren

Lawrence Berkeley National Laboratory
asalmgren@lbl.gov

Andrew Siegel
Argonne National Laboratory
siegela@uchicago.edu

Jean-Luc Vay
Lawrence Berkeley National Laboratory
Berkeley, CA
JLVay@lbl.gov

## MS9

### Microbiome Science at the Exascale

The ExaBiome project has developed exascale tools for analyzing microbial data. This includes a new metagenome assembler and tools for protein clustering, binning, and comparative analysis. For example, the MetaHipMer assembler allows for co-assembling massive environmental data sets with tens of terabytes of input data requiring petabytes of memory to assemble. With several record-breaking assemblies, the team recently assembled nearly 100 terabytes of ocean microbiome data and human gut microbiome data on the Frontier exascale system. Compared to previous approaches, the result is higher quality assemblies with more complete genomes, and more novel and rare species. Similarly, the HipMCL tool for protein clustering has been used on an unprecedented 8 billion proteins and resulted in the discovery of new protein clusters. These applications do not match the traditional patterns of scientific simulations but involve unstructured communication and computation, irregular data structures such as graphs and hash tables, and they require novel approaches to distributed memory parallelization and the use of GPU nodes.

Katherine Yelick
Lawrence Berkeley National Laboratory
University of California at Berkeley
yelick@eecs.berkeley.edu

<u>Leonid Oliker</u>, Steven Hofmeyr, Rob Egan, Aydin Buluc
Lawrence Berkeley National Laboratory
loliker@lbl.gov, shofmeyr@lbl.gov, rsegan@lbl.gov,
abuluc@lbl.gov

Oguz Selvitopi
Lawrence Berkeley Lab
roselvitopi@lbl.gov

Robert Riley, Richard Lettich
Lawrence Berkeley National Laboratory
rwriley@lbl.gov, rlettich@lbl.gov

Muaaz Awan
National Energy Research Scientific Computing Center
Lawrence Berkeley National Laboratory

mgawan@lbl.gov

## MS9

### Experiences from the ExaWind Project

The ExaWind project has developed multiple applications which are coupled together to perform blade-resolved wind farm simulations on the latest GPU-enabled supercomputers in the United States. ExaWind uses the AMR-Wind application for performing structured grid atmospheric flow between wind energy capturing turbines while Nalu-Wind uses unstructured meshes to resolve flow over each turbine. These applications are coupled through overset methodology, utilizing the TIOGA application to perform hole cuts and connectivity, while yet another application, Open-FAST, enables fluid structure interaction by translating loads on the turbine into mesh deformations. This talk will give an overview of the ExaWind project with specific detail on the successes and challenges experienced in the seven year project, as well as lessons learned.

Jon Rood, Michael Sprague
National Renewable Energy Laboratory
jon.rood@nrel.gov, michael.a.sprague@nrel.gov

## MS9

### Development of An Exascale Subsurface Simulator of Coupled Flow, Transport, Reactions Mechanics

Exascale computing has opened opportunities for modeling subsurface processes at scales that were previously not possible. Here we present an overview of the development of the Exascale Subsurface Simulator of Coupled Flow, Transport, Reactions Mechanics. We will discuss the challenges associated with moving the two component software codes (LBLs Chombo-Crunchflow and LLNLs GEOS.) onto exascale hardware. Strategies for performance portability will be presented, along with an overview of challenges faced by the project on the path to exascale. Finally, scaling studies for the project challenge problem will be presented with a discussion of continuing efforts to improve performance and scalability on the ORNL/Frontier exascale machine.

Randolph R. Settgast
Lawerence Livermore National Laboratory
settgast1@llnl.gov

David Trebotich
Lawrence Berkeley National Laboratory
dptrebotich@lbl.gov

Victor Magri, William Tobin
Lawrence Livermore National Laboratory
paludettomag1@llnl.gov, tobin6@llnl.gov

Terry J. Ligocki
Lawrence Berkeley Laboratory
TJLigocki@lbl.gov

## MS10

### Parallel Randomized Tucker Decomposition Algo-

### rithms

We propose to accelerate Tucker tensor decomposition algorithms by using randomization and parallelization. We present two algorithms that scale to large data and many processors, significantly reduce both computation and communication cost compared to previous deterministic and randomized approaches, and obtain nearly the same approximation errors. The key idea in our algorithms is to perform randomized sketches with Kronecker-structured random matrices, which reduces computation compared to unstructured matrices and can be implemented using a fundamental tensor computational kernel. We provide probabilistic error analysis of our algorithms and implement a new parallel algorithm for the structured randomized sketch. Our experimental results demonstrate that our combination of randomization and parallelization achieves accurate Tucker decompositions much faster than alternative approaches. We observe up to a 16X speedup over the fastest deterministic parallel implementation on 3D simulation data.

Grey Ballard
Wake Forest University
ballard@wfu.edu

## MS10

### Distributed Computation of Persistent (co)homology

Persistent homology captures the shape of input data at different scales. Typical inputs are point clouds or scalar functions on graphs/grids. Computation of persistent homology is essentially a special form of Gaussian elimination applied to a boundary matrix of a simplicial complex. The boundary matrix is very sparse in the beginning, but some of the columns become denser during the computation. Morozov and Lewis proposed a distributed algorithm that was based on a so-called blowup construction. However, duality (computing cohomology instead of homology) allows us to eliminate the combinatorially expensive blowup construction.

Arnur Nigmetov
Lawrence Berkeley National Laboratory
anigmetov@lbl.gov

## MS10

### Parallel Algorithms for Computing Electron Correlation Using Sparse Local Mp2 Method

Second order MllerPlesset perturbation theory (MP2) method is often used to compute electron correlation energy in computational chemistry. In this talk, we focus on using both distributed memory and shared memory parallel programming methods for the local MP2 method. Local MP2 allows us to transform the energy computation into solving a special sparse linear system. We design a sparsity-enforced Krylov method to solve the linear system iteratively, and allocate processors according to the sparse structures. We show strong scaling analysis on real chemical objects to indicate the efficiency and effectiveness

of our parallel algorithms.

Tianyi Shi
Lawrence Berkeley National Laboratory
tianyishi@lbl.gov

## MS10

### Riemannian Optimization for Tucker Tensor Completion

Tensor completion is a technique to fill in missing entries in multi-dimensional data using decomposition methods like Tucker, CP, or Tensor-train. Tucker tensor completion is used in applications when more accuracy is required, and the order of the tensor is not too large. Recent works have proposed Riemannian optimization algorithms for this problem and have shown that these algorithms perform better than traditional optimization methods like alternating minimization. We analyze the computations involved in Riemannian optimization algorithms and improve upon the asymptotic computational complexity in the existing algorithms. We leverage the sparse tensor contraction primitives to implement these algorithms and show that, with this improvement, Riemannian optimization can outperform alternating minimization in distributed-memory setting for practical applications.

Navjot Singh
University of Illinois at Urbana Champaign
navjot2@illinois.edu

## MS11

### Uncertainty in Uncertainty and Rockafellian Relaxation

A critical aspect of PDE constrained optimization is to account for uncertainty in the underlying physical models, for example in model coefficients, boundary conditions, and initial data. Uncertainty in physical systems is modeled with random variables, however, in practice there may be some nontrivial ambiguity in the underlying probability distribution from which they are sampled. As stochastic optimal control problems are known to be ill-conditioned to perturbations in the sampling distribution, we describe an analytic framework that is better conditioned to such meta-uncertainties and conclude with numerical examples.

Sean P. Carney
George Mason University
scarney6@gmu.edu

## MS11

### Derivative Tensor Network Compression via Symmetric Actions for Constructing High Order Taylor Series Surrogates

Digital twins for complex systems are often enabled by surrogates based on local linearization and low-rank approximation of the Jacobian. This technique neglects nonlinear behavior which may be essential to faithfully representing the system of interest. Taylor series surrogates capture important nonlinear behavior locally via higher order derivatives; however, the scale of high order derivative tensors poses several computational challenges. This work addresses these challenges via a derivative tensor network compression algorithm that relies solely on highly efficient and parallelizable symmetric derivative actions. The resulting method is applied to accelerate sampling from the posterior in a distributed parameter PDE-constrained inverse problem.

Blake Christierson
University of Texas at Austin
bechristierson@utexas.edu

Nicholas Alger
The University of Texas at Austin
Oden Institute
nalger225@gmail.com

Peng Chen
Georgia Institute of Technology
pchen402@gatech.edu

Omar Ghattas
University of Texas at Austin
omar@oden.utexas.edu

## MS11

### Local Point Spread Function Approximation of High-Rank Hessians in Inverse Problems Governed by PDEs

In this talk, we discuss the problem of estimating solutions of large-scale inverse problems governed by partial differential equation (PDE) based models. Solving inverse problems with expensive forward models and high-dimensional parameters is computationally challenging. One such challenge is that the Hessian, an object that is needed for a Newton-based solution of the inverse problem, is a matrix-free linear operator that requires the solution of two linear discretized PDEs in order to apply it to a vector. Here, we exploit the local sensitivity of model predictions to parameters which suggests that the Hessian has numerically low-rank off-diagonal blocks. We discuss an a priori means to estimate the supports of Hessian columns and how that knowledge can be leveraged for efficient batch sampling of its columns. This computational framework is ultimately utilized to estimate entries of the Hessian and compress it into a hierarchical H-matrix format. We apply such Hessian approximations as preconditioners in Newton linear system and ultimately reduce the total number of PDE solves required to solve two model inverse problems by the inexact Newton-CG method.

Tucker Hartland
Lawrence Livermore National Laboratory
hartland1@llnl.gov

Nicholas Alger
The University of Texas at Austin
Oden Institute
nalger225@gmail.com

Noemi Petra
University of California, Merced

npetra@ucmerced.edu

Omar Ghattas
University of Texas at Austin
omar@oden.utexas.edu

## MS11

### Fast and Scalable FFT-Based GPU-Accelerated Algorithms for Hessian-Vector Products Arising in Inverse Problems Governed by Autonomous Dynamical Systems

Hessian-based algorithms for the solution to inverse problems typically require many actions of the Hessian matrix on a vector (matvecs). A direct approach is often computationally intractable for problems with high-dimensional parameter fields or expensive-to-evaluate forward models. For systems that exhibit shift-invariance (e.g. autonomous systems) structure in their discretized form, the discretized linear parameter-to-observable (p2o) maps are block Toeplitz matrices. Moreover, considering causality for time-invariant systems the p2o map and its adjoint are lower- and upper-triangular block Toeplitz, respectively. By exploiting this structure, Hessian matrices for these types of systems can be compactly represented and Hessian matvecs can be efficiently computed through scalable multi-GPU FFT matvecs. Exploiting the triangular block Toeplitz structure yields memory savings proportional to the number of time steps $N_t$ and a computational speedup of $O(N_t/\log(N_t))$. We develop a multi-GPU FFT matvec code for Hessians corresponding to block Toeplitz p2o matrices utilizing the cuFFT and NCCL libraries. Our implementation achieves 75-90% of the maximum memory bandwidth on NVIDIA A100 80GB GPUs for all custom GPU kernels — which correspond to memory-bound operations. We also show strong and weak multi-GPU scaling on the Frontera RTX nodes with up to 81 GPUs.

Sreeram R. Venkat
University of Texas Austin
201 E 24th St, Austin, TX 78712
srvenkat@utexas.edu

Milinda Fernando, Stefan Henneking
Oden Institute, UT Austin
milinda@oden.utexas.edu, stefan@oden.utexas.edu

Omar Ghattas
University of Texas at Austin
omar@oden.utexas.edu

## MS12

### Parallel Exponential Polynomial Integrators

Exponential integrators are a class of time integration methods for efficiently solving stiff ordinary differential equations. A key characteristic of these methods is that they treat a linear component exactly while approximating all remaining terms explicitly. In this talk, I will present a new class of exponential integrators that are parallelizable both across the method and across the steps. The methods are based on the polynomial time integration framework and can be constructed at arbitrary orders of accuracy. I will discuss the stability of these methods, and present several numerical experiments that highlight the advantages of these integrators over existing parallel-in-time methods.

Tommaso Buvoli
University of California, Merced
tbuvoli@tulane.edu

## MS12

### Direct Time-Parallel Methods Based on Spectral Deferred Correction

During this talk, I will present one approach currently developed to perform PinT time-integration for numerical weather prediction. It is based on Spectral Deferred Correction (SDC), a time-integration method that iteratively computes the stages of a fully implicit collocation method using a preconditioned iteration. SDC allows to generate a variety of methods with arbitrary order of accuracy. While there are several parameters that can be used to optimize a SDC method, the main one is the choice of preconditionner. In particular, one can build diagonal SDC preconditioners, either to improve convergence speed or numerical stability for larger time-step size. One important aspect of those diagonal preconditioners is that they allow computations for SDC iterations to be performed in parallel in time. I will present some recent results on building optimized diagonal preconditionner for a split implicit-explicit formulation of SDC and show their application to test problems based on the shallow water equations.

Philip Freese, Thibaut Lunet
Hamburg University of Technology
philip.freese@tuhh.de, thibaut.lunet@tuhh.de

Daniel Ruprecht
Institute of Mathematics
Hamburg University of Technology
ruprecht@tuhh.de

## MS12

### Estimating Discrete Events for Spectral Deferred Corrections

Differential equations arise naturally in many fields to describe physical systems or processes. Discontinuities are also belong to such a system. In simulation, such discrete events can make the computation of a solution to a very difficult task. Even when the occurence of a discontinuity depends on the system dynamic itself, numerical methods may have problems calculating an accurate solution. Often they occur in fields such as science and engineering. In engineering, converters use high-frequency switching to compare after fixed number of time steps if the actual voltage is equal to a target voltage to control the output. In science, a gas-liquid model describes the output of a tube. The liquid comes out of the tube if the level of the liquid exceeds the dip tube, otherwise gas will leak out. In this presentation, the switch estimator applied to the method of spectral deferred corrections (SDC) is presented together with different examples to show its performance. The switch estimator predicts the time point of the discrete event us-

ing interpolation and root finding techniques, and the time step is adapted. Restarting SDC with the new time step the method is able to resolve the discontinuity more accurate.

Matthias Bolten
University of Wuppertal
bolten@math.uni-wuppertal.de

Kyrill Ho
Department of Mathematics and Computer Science
University of Cologne
kyrill.ho@uni-koeln.de

Robert Speck
Juelich Supercomputing Centre
Forschungszentrum Juelich GmbH
r.speck@fz-juelich.de

Lisa Wimmer
University of Wuppertal
wimmer@math.uni-wuppertal.de

Junjie Zhang
Institut für Energie- und Klimaforschung
Forschungszentrum Juelich GmbH
ju.zhang@fz-juelich.de

**MS13**

**Evolving Task-Based Parallel Programming to Address a Wider Range of Applications**

Task-based parallel programming models have been successful over the last fifteen years in helping port a range of applications to heterogeneous HPC platforms, thanks to their cooperation with dynamic runtime systems. Some applications have proven more challenging for task-based models, however, for a variety of reasons such as their inherent parallelism being more strenuous to extract or too overwhelming to manage. This talk will discuss some ideas we are exploring to enrich the programming model of our StarPU task-based runtime system to target new classes of applications.

Olivier Aumage
INRIA Bordeaux, France
olivier.aumage@inria.fr

**MS13**

**R&D in the Open—Tales from the LLVM Community**

The LLVM compiler framework is the backbone of all modern HPC vendor compilers and further used by most major IT companies. LLVM also underpins compilers for languages like Julia, Rust, and Switft, which makes it a very versatile target for research and development of novel compiler technology. In this talk we will provide insights about the R&D process in the LLVM community, especially with the HPC related work performed at DOE national laboratories. We will describe available resources, best practices, and show success stories. The latter will introduce the audience to various HPC/GPU related developments in LLVM which will have direct utility to them, including: a GPU libc and libm library, advanced performance portability for kernel languages and improved debugging and tuning support.

Johannes Doerfert
Lawrence Livermore National Laboratory
doerfert1@llnl.gov

**MS13**

**Ai Based Auto-Tuning for Iris Heterogeneous Runtime System with Matris Math Library Usecase**

IRIS is an Intelligent heterogeneous runtime system, enables task based programming model for applications to create tasks and run the kernels on heterogeneous compute units for high portability and performance. It abstracts the heterogeneous device memories management and enables automatic data movement with wide range of task schedulers to explore. MatRIS is a heterogeneous math kernel library developed using IRIS runtime and provides tiled algorithms for basic linear algebra. Though IRIS runtime and MatRIS provides several knobs such as tiling algorithm, tile size, number of heterogeneous compute units, type of compute units, etc., scheduling algorithm, etc. for each math kernel and for each input size, it is extremely challenging to select the right choice of knobs for best performance. We propose an generative AI based auto-tuning for the selection of these knobs for different problem input size of different math kernels. This talk covers the details of how AI algorithms help IRIS heterogeneous runtime and math kernels to identify the right set of knobs for different heterogeneous systems and provide some results.

Narasinga Rao Miniskar
Oak Ridge National Laboratory
miniskarnr@ornl.gov

**MS13**

**JACC.jl: A Julia CPU/GPU Portable Functional Layer**

We present JACC.jl, a functional-oriented portable programming model for the Julia language. JACC.jl provides a unified, lightweight front end across different back ends available in Julia (Base.Threads, CUDA.jl, and AMDGPU.jl) so that the same Julia code can run on different CPU and GPU targets. We evaluated the performance of JACC.jl for common HPC kernels (e.g., AXPY, conjugate gradient algorithm, Lattice-Boltzmann method) on nodes of the Summit (#5 TOP500) and Frontier (#1 TOP500) supercomputers at Oak Ridge National Laboratory. Overall, we show that the proposed programming model leverages newly introduced features in Julia v1.9 for optional dependencies while incurring a negligible performance overhead vs. Julia's vendor-specific solutions. We report expected GPU speedups over a CPU implementation with no extra cost to programmability for the kernel granularity. Hence, we attempt to leverage Julia as a high-productivity layer compiled via LLVM by providing a performance-portable solution.

Pedro Valero-Lara
Oak Ridge National Laboratory

valerolarap@ornl.gov

## MS14

### Optimizing Irregular Communication with Neighborhood Collectives and Locality-Aware Parallelism

Irregular communication often limits both the performance and scalability of parallel applications. Typically, the communication is implemented as point-to-point, and optimizations are integrated directly into the application, lacking portability. Optimization of point-to-point messages within MPI is difficult, as the interface only provides information on a piece of all communication. Persistent neighbor collectives offer an interface for optimizations within MPI. This paper presents methods for implementing existing optimizations for irregular communication within neighborhood collectives, analyzes the impact of replacing communication in existing codebases such as Hypre BoomerAMG with neighborhood collectives, and shows up to a 1.38x speedup on sparse matrix-vector multiplication through the use of our optimized neighbor collectives. The authors analyze three implementations of persistent neighborhood collectives for alltoallv: an unoptimized wrapper of standard point-to-point communication, and two locality-aware aggregating methods. The latter exposes an extended interface to perform additional optimization, and the authors show additional 7 percentage point speedup when using this non-standard interface.

Gerald Collom
University of New Mexico
Albuquerque NM, USA
geraldc@unm.edu

Amanda Bienz
University of New Mexico
bienz@unm.edu

Rui Peng Li
Lawrence Livermore National Laboratory
li50@llnl.gov

## MS14

### Performance Optimization of Sparse Factorization on GPUs

We present some results from the performance optimizations of the simplicial GPU-resident sparse LU and Cholesky factorizations available on in the Ginkgo software library. The limited amount of parallelism available in the dependency structure of a sparse factorization renders the common latency-hiding approach of GPUs ineffective, so we detail techniques for optimizing the handling of sparsity pattern lookups, load balancing between warps and between thread blocks, as well as special techniques for short rows. Finally, we show the local impact of these improvements based on a fine-grained profiling of the factorization using GPU clock registers.

Tobias Ribizel
Karlsruhe Institute of Technology

tobias.ribizel@kit.edu

## MS14

### Recent Experiences in Preconditioning for CFD on Hybrid Computing Architectures

In this talk we will overview recent efforts and outcomes in preconditioning linear systems arising from implicit Computational Fluid Dynamics (CFD) simulations on heterogenous computing architectures. Many problems in fluid dynamics require implicit methods to resolve disparate time scales for unsteady calculations as well as problems in which a steady state solution is required. With most of the compute power on modern supercomputers coming from Graphics Processing Units (GPUs), there is an imperative for developing algorithms and code for all classes of science problems which effectively utilize these architectures. Implicit computations, which heavily rely on linear solvers and preconditioners for performance, can be particularly difficult to map onto these new architectures. We will present performance results using several preconditioners for implicit CFD calculations with the matrices arising from applying Newtons method to the Navier-Stokes equations. Due to the poor conditioning of the linear systems, emphasis will be placed on using sparse direct linear solvers as a local subdomain preconditioner as part of a larger domain decomposition scheme. Both multithreaded CPU performance and GPU performance will be discussed.

Stuart Slattery
Computational Sciences and Engineering Division
Oak Ridge National Laboratory
slatterysr@ornl.gov

## MS14

### Re::Solve: Linear Solver Library Optimized for Solving Sequences of Linear Systems

The development of Re::Solve sparse linear solver library started as a part of Stochastic Grid Dynamics at Exascale (ExaSGD) project. The overarching goal was to address projects major technology gap at the time. The project required a fast sparse direct solver that is (natively) GPU-resident and at the same time, highly optimized, without unnecessary memory traffic nor repeated memory allocations and deallocations. A stable iterative refinement strategy was needed too, as the systems solved in ExaSGD were ill-conditioned by construction. During the project, it turned out that combining different codes and strategies (e.g., using LU decomposition from one solver library and following it with an alternative triangular solve and iterative refinement) is a winning approach and that a flexible library that facilitates this type of free mix-and-match solver style was needed. At the same time, the technology gap is not unique to ExaSGD, and affects many other applications whose overall performance heavily depends on the performance of the linear solver. Hence, Re::Solve was developed to be a versatile solver library that is flexible (e.g., same iterative solvers can be used for iterative refinement and direct solvers as preconditioners), designed to solve a sequence of linear systems without unnecessary recomputation and re-allocations, easy to integrate with applications, and capable of running on both AMD and

NVIDIA GPUs.

Katarzyna Swirydowicz
Pacific Northwest National Laboratory
kasia.swirydowicz@pnnl.gov

Slaven Peles
Oak Ridge National Laboratory
peless@ornl.gov

## MS15

### Dynamic Data Reduction in Motion: A Path to Efficient Sparse AI and Beyond

The ballooning energy consumption of AI urgently call for sustainable solutions. While domain-specific architectures offer some relief, they fall short of addressing the energy-draining issues related to data access and movement. This is especially true as AI models continue to scale and exhibit features like sparsity, which have traditionally been challenging in HPC applications. In this talk, we pivot away from traditional computing models to introduce a groundbreaking dynamic data reduction in motion paradigm. This novel approach dramatically reduces energy consumption by processing data dynamically as it traverses the system. It eliminates the need for energy-costly data movement, offering an anticipated 200-fold decrease in data movement-related energy consumption. Our system uniquely enables efficient parallel computation by optimizing data pathways, ensuring that data is processed as it moves from its original location to its destination. This dual focus on speed and efficiency sets a new benchmark for sustainable AI and large-scale scientific computing workloads that share features like sparsity. Preliminary results are promising. In the case of deep learning recommendation models, our approach outperforms existing methods in speed-up trends as system or batch sizes grow.

Bahar Asgari
University of Maryland
bahar@umd.edu

## MS15

### Accelerating Parallel and Distributed Systems and Applications with SmartNICs

The parallel and distributed computing community has been actively investigating methods to enhance system designs through the use of programmable SmartNICs, resulting in a significant boost in application efficiency. During this presentation, I will begin by outlining the opportunities and challenges associated with the design of distributed systems and applications utilizing modern SmartNIC architectures and technologies. Subsequently, I will delve into several case studies to illustrate our expertise in accelerating parallel and distributed systems and applications through the utilization of modern SmartNICs.

Xiaoyi Lu
University of California, Merced
xiaoyi.lu@ucmerced.edu

## MS15

### Lessons Learned from Accelerating Deep Learning Applications Using SmartnNICs

Many emerging deep learning models, such as recommendation systems, require large-scale multi-node environments for distributed training and inference. However, they often suffer from collective communication bottlenecks (e.g., alltoall), limiting their scalability. In this talk, I will discuss the lessons we've learned from designing a SmartNIC-based heterogeneous system. Through a software-hardware co-design approach, we have managed to overcome the communication bottleneck in distributed training, alleviate memory bandwidth pressure, and enhance computational efficiency.

Dingwen Tao
Indiana University
ditao@iu.edu

## MS15

### Intelligent Algorithms on Intelligent Networks-Experiences and Challenges Using Nvidias Blue-Field Technology

In this talk, we study HPC software which combines modern algorithmic ingredients such as dynamically adaptive mesh refinement and local time stepping with task-based parallelisation. In an ideal world, the software uses the tasks for lightweight load balancing, i.e. moves them between ranks. In reality, the code base suffers from high bandwidth latency and bandwidth constraints, and the MPI implementations struggle to cope with unexpected message arrival and MPI progression. We introduce our system architecture which deploys part of the task balancing and task migration to BlueFields, and we discuss to which degree our local BlueField software stack fits to these software ideas. First measurements for mini-app configurations demonstrate what the DPU system has to deliver in terms of latency and bandwidth to make our light-weight load balancing work.

Tobias Weinzierl
Durham University
tobias.weinzierl@durham.ac.uk

## MS16

### Accelerating Genome-Wide Association Study Using Mixed-Precision Fine-Grained Matrix Computations

Genome-wide association study (GWAS) represents a grand challenge due to the desire to incorporate millions of genetic markers over hundreds of thousands of patients to accurately model and predict genetic variations. Predictive models such as penalized regression is used to map human genetic variations to phenotypes. In particular, the Ridge Regression (RR) model involves solving systems of linear equations with large dense single nucleotide polymorphism (SNP) matrix obtained after forming the Gram matrix. The RR mapping estimates involve Level-3 BLAS opera-

tions, Cholesky factorization, and its corresponding solver for multiple right-hand sides. We develop a tile-based RR algorithm with mixed-precision techniques to accelerate successive computational phases. In particular, we employ an adaptive tile-centric mechanism to identify the precision arithmetic required for each tile. The resulting fine-grained computational tasks of RR can then be scheduled independently using the dynamic runtime system PaRSEC with asynchronous execution. We demonstrate how this controlled precision loss permits achieving the necessary accuracy for RR on a set of synthetic datasets. We further assess the numerical accuracy of the prediction on a real dataset from UK BioBank. We compare against popular RR approaches, such as NVIDIA Rapids and Regenie tools as well as the Deep Learning GenNet framework, and report accuracy and performance superiority on NVIDIA GPU hardware accelerators.

Rabab Omairy
King Abdullah University of Science and Technology
rabab.omairy@kaust.edu.sa

### MS17

### Accelerating Dense Density Matrix Computations on GPUs Using Chebyshev Expansions

Matrix diagonalization is often a bottleneck when computing the density matrix needed in electronic structure calculations. For modest matrix sizes (N= 4000 or less), performance of traditional dense diagonalization algorithms on modern GPUs is underwhelming compared to the peak performance of these devices. This motivates the exploration of alternative algorithms better suited to these types of architectures. Implementing a Chebyshev expansion algorithm whose number of required matrix multiplications scales with the square root of the number of terms in the expansion, our GPUs results show large speed ups compared to diagonalization for dense matrices. Additionally, we improve upon the original method by capitalizing on the inherent task concurrency in the algorithm using CUDA/HIP streams. This leads to a significant speed up over the serial-only approach for smaller (N=1000 or less) matrix sizes.

Jean-Luc Fattebert
Oak Ridge National Laboratory
fattebertj@ornl.gov

Joshua Finkelstein, Christian Negre
Los Alamos National Laboratory
jdf@lanl.gov, cnegre@lanl.gov

### MS17

### An Orbital-Wise Parallel Algorithm for Numerical Solutions of the Kohn-Sha Equation

Kohn-Sham density functional theory is one of the most successful approximate model for many-body Schrodinger equation. In this talk, an h-adaptive finite element framework for both Kohn-Sham equation and time-dependent Kohn-Sham equation will be described, based on which the effort towards the orbital-wise parallel will be introduced in detail, including the design of the preconditioner in solving generalised eigenvalue problem, a PDE-treecode

algorithm for Hartree potential calculation, a Jacobi-like update of orbitals in solving time-dependent equation, etc. Numerical experiments will demonstrate the effectiveness of our algorithm, and its potential in practical simulations.

Guanghui Hu
University of Macau
garyhu@um.edu.mo

### MS17

### Discontinuous Galerkin Hartree-Fock Calculations for Predicting Accurate Electronic Structures of Mesoscopic-Scale Metal-Semiconductor Junctions with Millions of Atoms

The evaluation of the exact Hartree-Fock exchange in hybrid density functional theory (DFT) is a crucial ingredient for accurately predicting electronic structures in molecules and solids. However, its application is currently limited to 5K atoms on leadership supercomputers due to its ultra-high computational complexity $O(N^4)$. Herein, we propose a new discontinuous Galerkin Hartree-Fock (DGHF) method for large-scale hybrid functional electronic structure calculations. We present a massively parallel DGHF implementation on exascale supercomputers to reduce the high computational scaling of constructing the HFX matrix from $O(N^4)$ to O(N). We showcase how DGHF can be used to predict accurate electronic structures of complex metal-semiconductor junctions with 2.5M atoms (17.2M electrons) using 35.9M cores on exascale Sunway supercomputer. This is the first time high-accuracy hybrid functional electronic structure calculations enable us to simulate next-generation electronic devices at mesoscopic scale (200 nm).

Xinming Qin, Junshi Chen, Wei Hu, Hong An, Jinlong Yang
University of Science and Technology of China
xmqin03@ustc.edu.cn,        juns@ustc.edu.cn,        whu-ustc@ustc.edu.cn, han@ustc.edu.cn, jlyang@ustc.edu.cn

### MS17

### Large-scale Density Functional Theory

Over the course of the past few decades, electronic structure calculations based on density functional theory (DFT) have become a cornerstone of materials research by virtue of the predictive power and fundamental insights they provide. The widespread use of the methodology can be attributed to its generality, simplicity, and high accuracy-to-cost ratio relative to other such ab initio approaches. However, while less expensive than wavefunction based methods, the solution of the DFT problem remains a formidable task. In this talk, the speaker will present recent developments in enabling large-scale DFT, including systems that contain over a million atoms and calculations that include many-body effects.

Phanish Suryanarayana
Georgia Institute of Technology

phanish.suryanarayana@ce.gatech.edu

## MS18

### A Performance Analysis of a High-Order 3D Immersed Interface Method Paired with a High-Order Wavelet-Based Adaptive Multiresolution Grid

We present a performance analysis of the high-order immersed interface method implemented within MURPHY, a C++/MPI software framework that provides 3D grid adaptation on distributed-memory machines. The implementation provides a high-order interface treatment on a block-based grid divided among multiple ranks, coupled with a high-order wavelet-based grid adaptation strategy. To compare the merits of both low and high order interface treatments, we evaluate the corresponding cost of geometry processing, bulk stencil operations, and altered boundary discretizations in terms of memory consumption, computation, and communication. Further, we evaluate the effect of low and high order grid adaptivity on these results. To conclude, we demonstrate the utility of the combined high-order boundary treatment and multiresolution techniques by presenting results from large-scale PDE simulations with complex geometries and fine-scale solution features on or near the immersed interface.

James Gabbard
Massachusetts Institute of Technology
jgabbard@mit.edu

Wim van Rees
MIT
wvanrees@mit.edu

## MS18

### Integrated Simulations of Compressible Multiphase Flows Using Diffuse Interface Methods and Data-Driven Spray Models

The overarching goal of the PSAAPIII INSIEME Center is the prediction of the reliability of laser-induced ignition of cryogenic propellants (liquid Oxygen and gaseous Methane) in a model rocket combustor in in-space conditions. This effort relies on simulations of this multi-physics problem involving multi-phase compressible fluid dynamics, thermodynamics, turbulent mixing, laser-induced ignition, and combustion at various levels of fidelity. Task-based programming in the HTR code (https://doi.org/10.1016/j.cpc.2020.107262), along with Legion as the software compiler and Exascale runtime environment, enables more seamless performance from next-generation heterogeneous supercomputers. In this talk, we will focus on the multiphase flow efforts within the Center, highlighting (1) the diffuse interface method used for capturing the evolution of the liquid jet and its primary atomization, (2) the data-driven model for atomization, which augments the interface-capturing method, and (3) the spray model which is of critical importance in predicting the spatiotemporal distribution of vapor mass fraction. By presenting simulation results and performance metrics, we will discuss the challenges involved in integrating these three components and our strategies for leveraging the unique physical conditions of the combustion chamber, the separation of scales involved in the two-phase jet, and the heterogeneous supercomputer architectures.

Shahab Mirjalili
Stanford University
ssmirjal@stanford.edu

## MS18

### Effective Gpu Utilization of a High-Order Embedded Boundary Method Through Brick Data Structures

We demonstrate an algorithm and data structures for embedded boundary schemes to effectively utilize GPU architectures. The high-order scheme is irregular, but well suited to take advantage of compute resources available on GPUs when data layouts are optimized. One of the significant challenges of high-order embedded boundary methods is the management of irregular data access patterns near boundaries; for this, we use a hash-based, brick data layout. This provides two significant advantages: irregularly shaped data is naturally supported with minimal excess memory, and data communication is efficient because no data packing is required. We show that this data structure is efficient on GPUs, even when being used for irregular domains. The approach supports complex embedded boundary schemes, which is demonstrated with implicit time integration for scalar diffusion. Comparisons of the algorithm run on both CPUs and GPUs are shown, and demonstrate portability. Performance analysis of the algorithm identifies both advantages of the design and directions for future improvement.

Nathaniel Overton-Katz
LBNL
NOvertonKatz@lbl.gov

Stephen Guzik
Colorado State University
stephen.guzik@colostate.edu

Hans Johansen
Lawrence Berkeley National Laboratory
Computational Research Division
hjohansen@lbl.gov

## MS18

### Scalability & Adaptivity: Achieving Conflicting Goals in a Heterogeneaous Computing Era

Partial differential equations (PDEs) are pervasive in engineering and science, and their numerical solutions are of paramount importance in understanding complex, natural, engineered, and societal systems. For large-scale PDE systems, adaptivity is essential to keep the memory footprint small, especially on modern architectures. At the same time scalability of our computational algorithms and codes is essential in order to be able to solve large problems. These two goals are however difficult to achieve on modern heterogeneous architectures. Adaptivity in particular leads to unstructured memory access leading to poor distributed-memory scalability. In this talk, I will present algorithms for large-scale adaptive refinement that achieve

high levels of spatial and temporal adaptivity without sacrificing scalability or efficiency. I will present a wide range of applications to demonstrate the efficacy of these methods, including thermal management of miniaturized integrated circuits, investigating the effect of platooning on transportation efficiency, analysis of COVID transmission in classrooms and simulations of jet atomization in the context of aerospace combustors.

Hari Sundar
School of Computing
University of Utah
hari@cs.utah.edu

## MS19

**Analysis of Dynamic Networks with Candy**

Updating properties of large dynamic networks present unique challenges. As the topology of the network changes, new algorithms are required to identify regions of change, to efficiently update the properties, as well as to maintain scalability. We are developing a software platform CANDY (Cyberinfrastructure for Accelerating Innovation in Network Dynamics) to analyze large dynamic graph efficiently. We will present the key algorithmic template of CANDY. Given a set of changed edges, which can be insertion or deletion, the template follows these two steps: The first step is graph sparsification, where each changed edge is processed in parallel to identify the affected vertices. The second step involves property update of affected vertices and maintaining correctness. The second step is an iterative process that uses parallel threads to operate on different affected vertices. We will demonstrate how this template can be applied to many different network properties including minimum weighted spanning tree, single source shortest paths, strongly connected components, page rank, and vertex coloring. We will discuss some of the future directions for updating dynamic networks, including the need for designing appropriate metrics and dynamic graph generators for accurate and standardized benchmarking.

Sanjukta Bhowmick
University of North Texas
sanjukta.bhowmick@unt.edu

Boyana Norris
University of Oregon
norris@cs.uoregon.edu

Sajal Das
Missouri Univeristy of Science and Technology
sdas@mst.edu

Sriram Srinivasan
University of Oregon
sriram882004@gmail.com

## MS19

**Engineering Practical Dynamic Graph Algorithms: the Parallel Batch-Dynamic Model**

In this talk, I will discuss the parallel batch-dynamic model for shared-memory multicore architectures. I will present several algorithms and implementations in this model for problems like k-core decomposition and triangle counting. Finally, I will discuss the improvements these algorithms obtain over state-of-the-art in practice.

Quanquan C. Liu
Yale University
quanquan.liu@yale.edu

Laxman Dhulipala
MIT
laxman@mit.edu

Jessica Shi
Massachusetts Institute of Technology
jeshi@mit.edu

Julian Shun, Shangdi Yu
MIT
jshun@mit.edu, shangdiy@mit.edu

## MS19

**Fused Breadth-First Probabilistic Traversals on Distributed GPU Systems**

Probabilistic breadth-first traversals (BPTs) are used in many network science and graph machine learning applications. We are motivated by the application of BPTs in stochastic diffusion-based graph problems such as influence maximization. These applications heavily rely on BPTs to implement a Monte-Carlo sampling step for their approximations. Given the large sampling complexity, stochasticity of the diffusion process, and the inherent irregularity in real-world graph topologies, efficiently parallelizing these BPTs remains significantly challenging. We present a new algorithm to fuse massive number of concurrently executing BPTs with random starts on the input graph. Our algorithm is designed to fuse BPTs by combining separate traversals into a unified frontier on distributed multi-GPU systems. To show the general applicability of the fused BPT technique, we have incorporated it into two state-of-the-art influence maximization parallel implementations (gIM and Ripples). Our experiments on up to 4K nodes of the OLCF Frontier supercomputer (32,768 GPUs and 196K CPU cores) show strong scaling behavior, and that fused BPTs can improve the performance of these implementations up to 34x (for gIM) and 360x (for Ripples).

Marco Minutoli
Pacific Northwest National Laboratory
marco.minutoli@pnnl.gov

Reece W Neff, Mostafa E Zarch
North Carolina State University
rwneff@ncsu.edu, meghbal@ncsu.edu

Mahantesh Halappanavar, Antonino Tumeo
Pacific Northwest National Laboratory
mahantesh.halappanavar@pnnl.gov,
Antonino.Tumeo@pnnl.gov

Ananth Kalyanaraman

Washington State University
ananth@wsu.edu

Michela Becchi
North Carolina State University
mbecchi@ncsu.edu

## MS19

### Learning Dynamic Temporal Graphs at Scale

Graph analytics can claim a large share of the credit for tackling many grand challenges – such as understanding the spread of pandemics, designing large-scale integrated circuits, rendering faithful cardiac simulations, and forecasting medium-range weather patterns. Dynamic and temporal information is often the key to extracting precise intelligence from real-world graphs. It is, however, also extremely difficult to build a computing system that can learn on dynamic temporal graphs. Notably, all three components of graph learning (i.e., inference, training, and sampling) face dramatic challenges: (i) inference on a temporal dynamic graph can only tolerate a sub-millisecond latency, whereas existing inference mechanisms are slow because they have to perform three time-consuming sub-tasks sequentially. (ii) It is often desirable to retrain the model daily for dynamic graphs, whereas training GNNs on a billion-edge graph would require a week. Not to mention the real-world trillion-edge graphs. (iii) Monte Carlo sampling, the fundamental primitive to enable inference and training on real-world trillion-edge graphs, faces mounting challenges on dynamic temporal graphs. In this talk, we will discuss three published papers addressing these issues, which permit learning dynamic temporal graphs at scale.

Hang Liu
Rutgers Univesity
hang.liu@rutgers.edu

Shiyang Chen, Santosh Pandey
Rutgers University
shiyang.chen@rutgers.edu, santosh.pandey@rutgers.edu

## MS20

### EXAGRAPH: Graph and Combinatorial Methods for Enabling Exascale Applications

Combinatorial algorithms in general and graph algorithms in particular play a critical enabling role in numerous scientific applications. However, the irregular memory access nature of these algorithms makes them one of the hardest algorithmic kernels to implement on parallel systems. With tens of billions of hardware threads and deep memory hierarchies, the exascale computing systems in particular pose extreme challenges in scaling graph algorithms. The codesign center on combinatorial algorithms, ExaGraph, was established to design and develop methods and techniques for efficient implementation of key combinatorial (graph) algorithms chosen from a diverse set of exascale applications. Algebraic and combinatorial methods have a complementary role in the advancement of computational science and engineering, including playing an enabling role on each other. In this presentation, we survey the algorithmic and software development activities performed under the auspices of ExaGraph from both a combinatorial and an algebraic perspective. We detail our recent efforts in porting the algorithms to manycore accelerator (GPU) architectures.

Mahantesh Halappanavar
Pacific Northwest National Laboratory
mahantesh.halappanavar@pnnl.gov

Erik G. Boman
Center for Computing Research
Sandia National Labs
egboman@sandia.gov

Aydin Buluc
Lawrence Berkeley National Laboratory
abuluc@lbl.gov

Alex Pothen
Purdue University
Department of Computer Science
apothen@purdue.edu

## MS20

### Co-Designing Amrex: Challenges and Successes

AMReX is a numerical software framework for developing massively parallel, block-structured adaptive mesh refinement applications that targets a variety of platforms, including exascale hardware, and was developed as part of the Exascale Computing Project's Block-Structured Adaptive Mesh Refinement Co-Design Center. AMReX forms the basis for the spatial and temporal discretization strategies for a large number of scientific and engineering simulation codes - including codes that are part of six ECP application development projects - spanning fields such as accelerator design, astrophysics, combustion, cosmology, wind energy, and multiphase flows. In this talk, I will highlight some of the successes and challenges encountered in co-designing AMReX along with our application partners to take advantage of current and upcoming supercomputing architectures. I will particularly focus on lessons learned in refactoring and redesigning codes to take advantage of hybrid CPU/GPU systems, with the goal that other framework developers and application teams can learn from our experience.

Andrew Myers
Lawrence Berkeley National Lab
atmyers@lbl.gov

Ann S. Almgren
Lawrence Berkeley National Laboratory
asalmgren@lbl.gov

John B. Bell
Lawrence Berkeley National Laboratory, U.S.
jbbell@lbl.gov

Weiqun Zhang
Lawrence Berkeley National Laboratory
Center for Computational Sciences and Engineering

weiqunzhang@lbl.gov

## MS20

### Helping Decarbonize the North American Power Grid with Exascale Computing

Decarbonizing the North American Power Grid will require increased penetration of renewable energy production, including wind and solar, which are tightly coupled to weather for local effects and climate for long-term effects. But weather and climate increase the uncertainty of producing power. As a result, safely and reliably delivering power in the grid of the future will depend on our ability to foresee many more possible scenarios than can be considered today. The current power grid is operated using optimization calculations that minimize the cost of generating power to meet an anticipated need, subject to safety and security constraints. This calculation is typically expected to complete within a 30 minute operational window in order to be actionable to power grid operators. An Exascale Computing Project application called ExaGO casts this as a nonlinear optimization problem with equality and inequality constraints. The FRONTIER Exascale computing architecture at Oak Ridge National Laboratory provides an enormous amount of compute power, making it possible to perform ExaGO calculations on power grid models the size of the Western North American Interconnect, and for large numbers of contingencies (possible failure states) and weather scenarios. This presentation will describe algorithmic and implementation lessons learned for performing these calculations on the GPU-dominated environment in todays Exascale systems.

Chris Oehmen
Pacific Northwest National Laboratory
christopher.oehmen@pnl.gov

## MS20

### Interleaving Artificial Intelligence and Simulation at Scale with Colmena

The increasing capabilities of Artificial Intelligence (AI) have opened many new routes for accelerating computational workflows. We created a computational-steering tool, Colmena, as part of the ExaLearn project to build applications that make full use of AI methods and can deploy across thousands of compute nodes. Colmena lets scientists express sophisticated policies for how to use a variety of AI techniques (e.g., generative models, optimal experimental design) to steer simulation workflows. The policies are enacted by workflow engines from the Exascale Computing Project and accelerated by a high-performance data transfer fabric, ProxyStore. Our talk will describe the challenges addressed in running Colmena applications that employ unique patterns of combining AI with simulation at large scale for applications including molecular design, protein generation, and training surrogate models for quantum chemistry. We expect that the lessons learned will be relevant to future applications that use AI for Science.

Logan Ward
Argonne National Laboratory, U.S.
lward@anl.gov

## MS21

### Distributed-Memory Randomized Algorithms for Sparse Tensor CP Decomposition

Randomized linear algebra methods can offer order-of-magnitude speedups over their non-randomized counterparts, but extending them to distributed-memory parallel machines poses unique challenges. We parallelize a randomized sampling algorithm for sparse Candecomp / PARAFAC decomposition, a popular technique to extract features from high-dimensional sparse data. Several prior works have analyzed computation and processor-to-processor communication costs in non-randomized CP decomposition, with highly-optimized software packages available. While randomized algorithms offer significant speedups, we demonstrate, both in theory and practice, that they suffer from higher relative communication costs and load imbalances in the distributed-memory setting. We adapt the data distribution and local sparse tensor storage format to mitigate these issues, allowing communication to scale with the number of random samples taken. Experiments on sparse tensors with billions of nonzero entries demonstrate significant speedups over state-of-the-art distributed-memory libraries for sparse CP decomposition.

Vivek Bharadwaj
University of California, Berkeley
vivek_bharadwaj@berkeley.edu

## MS21

### Accelerating Sparse Matrix Computations with Code Specialization

Sparse matrix methods are at the heart of many scientific computations and data analytics codes. Sparse matrix kernels often dominate the overall execution time of many simulations. Further, the indirection from indexing and looping over the nonzero elements of a sparse data structure often limits the optimization of such codes. In this talk, I will introduce code specialization strategies that transform computation patterns in sparse matrix methods for high-performance. Specifically, I will show how decoupling symbolic analysis from numerical manipulation will enable the automatic optimization of sparse codes.

Maryam Dehnavi
University of Toronto
mmehride@cs.toronto.edu

## MS21

### Distributed-Memory Algorithms For Graph Embedding and Visualization

Large scale data visualization and searching is a requirement in data science for many fields, such as bioinformatics, molecular science, etc. Existing data pipelines such as UMAP are shared memory implementations that could not handle billion-scale data nodes. The need for distributed data visualization pipeline implementation motivated us to work on that area. We have implemented a two-phase dis-

tributed memory pipeline for data visualization. Firstly, we created a K-nearest neighbor graph using sparse random projection trees. We used distributed multiple random projection trees to classify similar data nodes into buckets followed by a small-scale exact search within the same bucket of nodes and across trees. Secondly, we use this graph to embed data into low dimensional space using a novel distributed memory embedding algorithm using t-distribution and force models. The output of the embedding algorithm can be directly used to visualize the data.

Isuru Ranawaka
Indiana University
isjarana@iu.edu

## MS21

### Communication Optimal Direct Solvers for Planar Sparse Matrices

Solving a system of linear equations $Ax = b$ where the matrix $A$ is a large planar sparse matrix, as commonly found in 2D PDEs, is essential in numerous scientific applications. Parallel solving of these matrices is a critical substep, but achieving communication optimality in their direct solution has been a long-standing challenge. In this talk, we will explore the hurdles faced with planar sparse matrices, where existing methods, optimal for other classes of sparse matrices, falter. We will introduce innovative algorithms that overcome these challenges, achieving optimal communication in sparse LU factorization. Also, we will look over the communication requirements for related procedures such as sparse triangular solving and iterative refinement. If time permits, we'll compare direct and Krylov subspace-based iterative methods and discuss the optimal balance between communication, memory, and precision. By shedding light on these challenges and presenting potential solutions, this talk aims to offer new insights and open questions for optimal communication sparse matrix computation for researchers and practitioners in scientific computing.

Piyush Sao
Oak Ridge National Laboratory
saopk@ornl.gov

## MS22

### Speeding Up Bayesian Inverse Problems with Fused Simulations

Solving Bayesian inverse problems with Markov Chain Monte Carlo (MCMC) methods has the advantage that only forward models have to be evaluated, and no adjoint problems have to be solved. At the same time, a lot of these forward models have to be evaluated, which requires an efficient solution method. SeisSol offers so called "fused simulations', where several earthquake scenarios are simulated at once. By using memory and SIMD vectors more efficiently, this approach speeds up the computation time per simulation. The Generalized Metropolis Hastings algorithm (GMH) is a parallelized extension to well-established Metropolis Hastings sampling, which evaluates several samples at once and adds the best fitting subset to the Markov chain. Fused simulations and the GMH kernel accommodate each other perfectly. We will present the solution of a kinematic source inversion problem with the GMH kernel and fused simulations. Furthermore, we will show how the method can be extended to include dynamic earthquake sources.

Michael Bader
Technical University of Munich
bader@in.tum.de

Sebastian Wolf
Technical University Munich
s.wolf@tum.de

Vikas Kurapati
Technical University of Munich
vikas.kurapati@tum.de

## MS22

### Parallel Algorithms for Experimental Design in Semi-supervised Learning

We propose a parallel algorithm for pool-based active learning for multiclass classification using multinomial logistic regression. Using finite sample analysis, we prove that the Fisher Information Ratio (FIR) lower and upper bounds the excess risk. Based on our theoretical analysis, we propose a parallel algorithm that employs regret minimization to minimize the FIR. We show algorithmic scalability for a multiclass classification problem that can scale to billions of points and thousands of classes. Combined with an experimental design algorithm, the new scheme can be used to select scenarios for training for digital twin applications.

George Biros
University of Texas at Austin
Oden Institute
biros@oden.utexas.edu

Youguang Chen
University of Texas at Austin
youguang@utexas.edu

## MS22

### Advancing Earthquake Physics through HPC-Enabled Digital Twins

Geo-hazards and risks increase worldwide rapidly due to continuing urbanization, climate change, and high-risk critical distributed infrastructure. The longest modern instrumental records of earthquakes cover less than 100 years, while recurrence intervals of large earthquakes are hundreds of years or more. Increasingly dense observations and physics-based simulations empowered by supercomputing provide pathways for overcoming the lack of data and elucidating spatiotemporal patterns that extend our knowledge beyond sporadic case studies and average statistical laws - however, are typically challenging to integrate. Digital Twins are emerging in Solid Earth Science, allowing curiosity-driven science to test scientific hypotheses against observations over ranges of space-time scales not accessi-

ble for laboratory and field observations. I will present results of large-scale earthquake simulations that aim to clarify the physical processes governing large earthquakes, improve our forecasting ability, and enhance the general understanding of earthquakes and faults. A prime example are our HPC-empowered complex rupture dynamics and ground shaking simulations of the two large, multi-fault earthquakes of the 2023 Kahramanmaras, Turkey, earthquake doublet using early observations.

Alice Gabriel
UC San Diego
liese.gabriel@gmail.com

## MS22

### Inference of Rheological Parameters in 1-Km-Resolution Earth Mantle Models

Estimating parameters in mantle flow and plate tectonics from surface observations results in an optimization problem. The forward problem for mantle flow is governed by highly nonlinear, heterogeneous, and incompressible Stokes equations. Solving these governing equations is already a major challenge at extreme computing scales. Adding an outer loop for parameter estimation adds substantially to the solver challenges. The computational methods for the forward problem rely heavily on adaptive meshes for local 1-km-resolution of the globe. The methods further include inexact Newton-Krylov with a combination of "BFBT" and multigrid preconditioning for saddle point linear systems. Scalable parameter estimation is enabled by derived adjoint Stokes equations within a Newton's method. Uncertainties of parameters are revealed by local approximations at the MAP point by computing a Gauss-Newton Hessian. We show inference on cross-sectional models of the Pacific and the first global inference results for 1-km-resolving models.

Johann Rudi
Virginia Tech
jrudi@vt.edu

Jiashun Hu
Southern University of Science and Technology
hujs@sustech.edu.cn

Michael Gurnis
Seismological Laboratory
California Institute of Technology
gurnis@gps.caltech.edu

Georg Stadler
Courant Institute for Mathematical Sciences
New York University
stadler@cims.nyu.edu

## MS23

### Pitsbicg: Parallel in Time Stable Bi-Conjugate Gradient Algorithm

This work introduces a novel algorithm for the parallel-in-time (PinT) numerical simulation of time-dependent partial and ordinary differential equations. Departing from the established parareal in time algorithm, we present a robust alternative by algebraically formulating the PinT problem and leveraging an adapted Bi-Conjugate Gradient Stabilized method. Referred to as the Parallel in Time Stable Bi-Conjugate algorithm (PiTSBiCG), this method exhibits substantial potential for stabilizing parallel resolutions across diverse problems. In this presentation, we elucidate the mathematical foundation of PiTSBiCG and substantiate its superiority over the conventional parareal approach through compelling numerical evidence.

Mohamed K. Riahi
Khalifa University of Science and Technology
mohamed.riahi@ku.ac.ae

## MS23

### A New Relaxation Scheme for Parallel-in-Time Methods with Absolute Convergence

Multigrid-in-Time methods such as Multigrid Reduction in Time (MGRIT) utilize multiple time-grids of varying resolution, combining a parallel relaxation method on fine-grids with a cheap, sequential solve on the coarsest grid, in order to eliminate the many sequential time steps that are normally required to solve time-dependent systems. While classical MGRIT has demonstrated optimal scaling for parabolic problems and recent work [H. De Sterck, et al.] has shown promise for hyperbolic systems, chaotic systems remain difficult to solve. This is because these systems are linearly unstable, so errors grow exponentially fast in time such that MGRIT convergence is much slower at the end of the time domain than the beginning, eventually stalling over longer time horizons. While recent work by the presenters has partially alleviated this problem by modifying the coarse-grid equation, there is still a fundamental time limit. Here we present a novel relaxation scheme, called Least Squares Relaxation (LSR), which relaxes the residual even for linearly unstable problems, in contrast with the classical FCF-relaxation used by MGRIT. Using LSR, the initial condition is partially relaxed so that information can flow symmetrically forward and backward in time, in principle allowing MGRIT to converge for chaotic problems on arbitrarily long time domains. We explore the properties of LSR and demonstrate its efficacy on the classical Lorenz system, as well as discuss considerations for PDEs.

David Vargas
University of New Mexico
dvargas2@unm.edu

Robert D. Falgout, Stefanie Guenther
Lawrence Livermore National Laboratory
falgout2@llnl.gov, guenther5@llnl.gov

Jacob B. Schroder
Department of Mathematics and Statistics
University of New Mexico
jbschroder@unm.edu

## MS23

### Time-Parallel Multigrid Preconditioning for KKT

### Systems Arising in Constrained Optimization

Optimal control problems governed by time-dependent partial differential equations (PDEs) lead to large-scale optimization problems. The cost of numerically solving these problems is proportional to the size of the discrete time domain. One approach to solve such problems is the use of algorithms based on the augmented systems. In the context of time-dependent problems, an augmented system is the KKT system for a convex linear-quadratic problem of type with a specific form of the objective function. However, in most numerical approaches the time discretization serializes the solution process and introduces a bottleneck. In the strong scaling limit, this bottleneck cannot be overcome by additional parallelization in space. To accelerate the solution of linear-quadratic optimal control problems governed by PDEs, we propose a time-domain decomposition approach that introduces time parallelism into the optimization algorithm. This is achieved by introducing auxiliary state variables for each time interval and impose time-continuity constraint. Additionally, auxiliary variables are incorporated into the objective function. We demonstrated the potential effectiveness of our scheme as a time-parallel preconditioner for linear systems arising in inexact SQP.

Radoslav G. Vuchkov
Sandia National Laboratories
rgvuchk@sandia.gov

Eric C. Cyr
Computational Mathematics Department
Sandia National Laboratotories
eccyr@sandia.gov

### MS23

### Parallel Performance of Block Epsilon-circulant Preconditioner for Time-dependent PDEs

This work considers the parallel performance of block epsilon-circulant (BEC) preconditioner for an all-at-once linear system arising from time-dependent PDEs. This preconditioner is a block circulant preconditioner in time and introduces the weight parameter epsilon at the top-right blocks. The diagonalization property of the preconditioned system leads to rigorous convergence analysis even if the original system is non-diagonalizable, and the preconditioned solver achieves convergence independent of the problem size. This work discusses the parallel-in-time performance of the BEC preconditioned GMRES solver. One of the main components is a one-dimensional parallel FFT for vectors of time-step size. We use three types of FFTs implemented using the libraries FFTW and FFTE: the straightforward method, the six-step FFT, and the redistributed FFT. Numerical experiments show parallel results for diffusion and convection-diffusion problems compared to a time-sequential solver and a multigrid-based parallel-in-time solver.

Ryo Yoda
Bergische Universität Wuppertal
ryoda@uni-wuppertal.de

Matthias Bolten
University of Wuppertal
bolten@math.uni-wuppertal.de

### MS24

### Portable Performance in General Purpose Exascale Calculations Challenges and Solutions?

The challenge of both porting and obtaining consistent performance from new heterogeneous exascale-and-beyond architectures is addressed by an analysis of the components of the architectures and the performance imbalances that they carry. This leads to consideration of coding upon software portability abstractions and hiding delays that damage scalability. The benefits of this approach are shown both through high performance GPU solutions using these approaches and through performance comparisons across different GPUs using a single code base. A look forward to future architectures is taken. This talk is based on work by about a dozen different collaborators in Utah, NIST and DOE.

Martin Berzins
Scientific Computing and Imaging Institute
University of Utah
mb@sci.utah.edu

### MS24

### Distributed Parallel Productivity Computing with Kokkos

Device-initiated Communication with RMA is a promising approach to support distributed parallel programming with modern data abstractions such as std::mdspan. In this talk, we give an introduction to Kokkos and its distributed memory support through Kokkos Views and Kokkos Remote Spaces. We also show the required API to implement data mapping, block transfers, and synchronization and showcase the envisioned use case with relevant applications on modern GPUs.

Jan Ciesko
Sandia National Laboratory
jciesko@sandia.gov

### MS24

### Clacc: Openacc Support for C/C++ in Clang and Llvm

Clacc has developed an open-source OpenACC compiler, runtime, and profiling support for C/C++ by extending Clang and LLVM under the Exascale Computing Project (ECP). A key Clacc design decision is to translate OpenACC to OpenMP in order to leverage the OpenMP offloading support that is actively being developed for Clang and LLVM. A benefit of this design is support for two compilation modes: traditional compilation mode translates OpenACC source to an executable, and source-to-source mode translates OpenACC source to OpenMP source. The purpose of this talk is to present the current status of the Clacc project, including recent support for OpenACC in C++, support for Kokkos's OpenACC backend, and Clacc's deployment on ORNL's Frontier, where Clacc is currently the only OpenACC implementation for C/C++.

Joel E. Denny
Oak Ridge National Laboratory
dennyje@ornl.gov

## MS24

### Omni Source-to-Source Compiler Infrastructure for Today's and Future Programming Models

Omni Compiler is an infrastructure for source-to-source transformation to design source-to-source compilers. It enables directive-based extension for C and Fortran95. We have been working on XcalableMP and XcalableACC, a directive-based language extension of Fortran95 and C for high-performance distributed memory parallel systems with accelerators, using Omni compiler. Recently, we have been working on OpenMP targets and OpenACC for FP-GAs. In this talk, the overview of Omni Compiler and a perspective for future programming models by directive extensions using Omni compiler will be presented.

Mitsuhisa Sato
RIKEN Advanced Institute for Computational Science
msato@riken.jp

## MS25

### Nvidia Accelerated Direct Sparse Solver

For many years, it was thought that direct sparse solver algorithms are not a fit for acceleration. To demonstrate the opportunity of acceleration we implemented a new Direct Sparse Solver Library for Nvidia GPU to based on the classical multifrontal approach. Our library supports typical DSS functionality as reordering/symmetric and non-symmetric factorization, transpose and non-transpose and many right-hand sides. Each feature has given about a 10x performance improvement on customer applications compared to native or unaccelerated solutions. We highlight global pivoting on a non-symmetric matrix, allowing us to solve systems that were not solved before. Performance comparison with major Direct Sparse Solver solutions for CPU and GPU presented.

Anton Anders, Alexander Kalinkin
NVIDIA
antona@nvidia.com, akalinkin@nvidia.com

## MS25

### Solving under-Constrained Hyperelasticity Without the Null Space

When hyperelastic structures are under-constrained, the system matrix has a nontrivial null space. To solve these singular linear systems with preconditioned iterative methods, one needs to exactly specify the null space. However, there are cases, such as contact and slip/symmetry boundary conditions, where identifying the null space analytically is error-prone and cumbersome. We propose an algorithm for which the user only needs to specify the rigid body modes (near-null space, which is already needed for algebraic multigrid) and the null space is computed automatically as the body gets deformed. We have implemented

this method in PETSc, the Portable Extensible Toolkit for Scientific computing. In this talk, we explore the effectiveness of this technique with p-multigrid using Ratel, a new solid mechanics library based on libCEED and PETSc.

Leila Ghaffari, Jeremy Thompson
University of Colorado Boulder
leila.ghaffari@colorado.edu,
jeremy.thompson@colorado.edu

Zachary R. Atkins
University of Kansas
zatkins@ku.edu

Evan Gassiot
University of Colorado
evan.gassiot@colorado.edu

Jed Brown
University of Colorado Boulder
jed@jedbrown.org

## MS25

### Stopping Criteria for the Conjugate Gradient Algorithm in High-Order Finite Element Methods

The finite element discretization of partial differential equations results in solving large-scale linear systems using iterative methods. This introduces errors from both the discretization and the inexact linear solve. Achieving an optimal balance between these errors requires a stopping criterion for the iterative method that includes an appropriate error indicator for the finite element approximation. However, the commonly used criterion based on the relative residual norm with a preset tolerance disregards the discretization error, often oversolving the linear system and leading to wasteful iterations. In this talk, we present a stopping criterion for high-order finite element discretization. Our criterion effectively identifies the optimal stopping point, regardless of mesh size, polynomial degree, or mesh shape regularity. We discuss our GPU implementation of the criterion and illustrate its efficacy and performance within libParanumal, an experimental suite of finite element flow solvers for heterogeneous GPU/CPU systems.

Yichen Guo, Eric De Sturler, Tim Warburton
Virginia Tech
ycguo@vt.edu, sturler@vt.edu, tcew@vt.edu

## MS25

### Distributed Sparse Solvers on AMD GPUs

Sparse direct solvers are needed widely in scientific computing, including e.g. discretized PDEs (CFD and FEA), bioinformatics, and machine learning. A computationally expensive step in the solution of a sparse set of linear equations is decomposition of the coefficient matrix into a product of simpler matrices; for a symmetric positive definite matrix, such a method is the Cholesky factorization. While GPUs are used widely to accelerate parallelizable problems such as matrix factorization, obtaining good GPU utilization with Cholesky factorization acting on a sparse matrix

has been a historically challenging problem, requiring algorithmic innovation to achieve performance. In this talk, we discuss optimizations to sparse Cholesky factorization as implemented in AMD's rocSOLVER routines. Furthermore, we present benchmarking results on AMD hardware, and performance of the sparse Cholesky factorization insitu from the ECP Project ExaSGD, an exascale application that enables the near real-time response to induced stresses on the national energy grid.

Michael Rowan
AMD
michael.rowan@amd.com

## MS26

### Accelerating I/O Services Using Smartnics (tentative)

SmartNICs are more than just additional compute and memory resources attached to a compute node. They can also be self-hosted solutions that can work in conjunction with existing host resources to help manage inter-node resources and coordination. As the SmartNIC is not typically powerful enough to mix into the application resources and use for compute cycles for the primary application. This talk will address the opportunities of regarding a SmartNIC as a "helper' system inside the node that has inter-node coordination capabilities. Using this frame of reference for SmartNICs, we can develop several "killer app' solutions for them and motivate their further integration into the data center of tomorrow.

Ryan Grant
Queen's University
ryan.grant@queensu.ca

## MS26

### SmartNIC-Accelerated Molecular Dynamics and Algebraic Multigrid Algorithms

In the same way that GPUs led to new (or revived) algorithmic methods to exploit them, we believe the same will be necessary for smartNICs. This talk presents this position in the context of two proxy applications, one from molecular dynamics and the other from an algebraic multigrid solver for electromagnetic analysis problems.

Sara Karamati, Jeffery Young, Rich Vuduc
Georgia Institute of Technology
s.karamati@gatech.edu, jyoung9@gatech.edu, richie@cc.gatech.edu

## MS26

### Offloading Data Management Services to Smartnics (tentative)

Data management and analysis is a critical part of large-scale science and engineering applications. These applications produce vast amounts of data that require complex analysis. Leveraging powerful, existing tools requires data pre-processing to match tool requirements. However, these data processing and transformation tasks consume resources that would otherwise be available to advance the application's computation. Smart Network Interface Cards (SmartNICs) provide valuable computational resources to scientific applications. Although the these devices are typically much less powerful than the associated host, SmartNICs have the potential to efficiently support offloaded computation that may not require the full compute power of the host. In this presentation, we consider the potential of offloading data management and analysis tasks to SmartNICs. Specifically, we consider how to efficiently move application output data from host to device memory. We also discuss a case study based on offloading Apache Arrow task to NVIDIA BlueField DPUs.

Scott Levy
Sandia National Laboratories
sllevy@sandia.gov

## MS26

### ODOS: Supporting DPU Offloading with OpenMP Directives

Data processing units (DPUs) as network co-processors are an emerging trend in our community, with plenty of opportunities yet to be explored. These have been generally used as domain-specific accelerators transparent to application developers; In the HPC field, DPUs have been used as MPI accelerators, but also to offload some tasks from the general-purpose processor. However, the latter required application developers to deploy MPI ranks in the DPUs, as if they were remote (weak) compute nodes, hence considerably hindering programmability. The wide adoption of OpenMP as the threading model in the HPC arena, along with that of GPU accelerators, is making OpenMP offloading to GPUs a wide trend for HPC applications. We will present ODOS, our brand-new OpenMP offloading support for network co-processor DPUs. We present our design in LLVM to support OpenMP standard offloading semantics and discuss the programming productivity advantages with respect to the existing MPI-based programming model. We also provide the corresponding performance analysis demonstrating competitive results in comparison with the MPI baseline.

Antonio J. Peña, Sergio Iserte, Muhammad Usman
Barcelona Supercomputing Center (BSC)
antonio.pena@bsc.es, sergio.iserte@bsc.es, muhammad.usman@bsc.es

## MS27

### Parallel Sparse Matrix Kernels for Graph Machine Learning

Sparse matrix operations such as sparse-dense and sampled dense-dense matrix multiplications play crucial roles in algorithms for graph embedding, graph neural networks and graph visualization. We develop parallel sparse matrix algorithms to accelerate graph machine learning. Our kernels are autotuned based on hardware features and matrix sparsity with an aim to perform better for different datasets on heterogeneous computing platforms. We plugged our sparse kernels with PyTorch and demonstrated significant speedups over kernels offered by PyTorch Geometric.

Ariful Azad

Indiana University
azad@iu.edu

## MS27

**Computational Study of Matrix Decomposition Methods for a Compression of a Transformer Neural Network Architecture**

The field of Natural Language Processing has made significant progress with the development of large language models based on Transformer architecture. Nevertheless, these models share a common challenge of expanding scale, presenting a formidable obstacle to model training. This scalability issue poses a bottleneck for scientific progress, impacting not only large-scale industries but also smaller research teams lacking comparable training resources. By optimising the training process, we use matrix (SVD) and tensor (TTM) decomposition techniques to represent some layers in the pre-trained Transformer models. To avoid drawbacks in the model performance, we align the compression objective and the task objective by injecting Fisher Information into the decomposition algorithms. We considered the time and memory required for forward-backwards signal propagation within the decomposed layer, as well as the quality obtained by the compressed model on downstream tasks. Measurements were carried out at various compression levels as well as different forms of the TTM decomposition.

Viktoriia Chekalina
Skoltech
vi.chekalina@skoltech.ru

## MS27

**A Task-Based Approach for a Parallel Distributed Training of Large Neural Networks**

Many classical problems are reduced to certain operations with systems of linear equations. Results with improved accuracy are usually achieved by straightforward increase of an amount of these linear equations, e.g., up to millions or even billions of rows and columns. Lots of HPC software rely on this fact. Unlikely, large neural networks are based on long chains of multiplications by relatively small matrices, e.g. 1600-by-6400 for the GPT2-XL model. This is one of the main reasons why performance of training stage barely surpasses 50-percent efficiency, meaning computing hardware is idle nearly half of the time. Splitting model weights into chunks to improve parallelism might even reduce efficiency, as small chunks of data inherit low arithmetic intensity. Therefore implementing task-based parallelism for training and inference of large neural networks requires overhaul of a task-based approach, especially task-scheduling algorithms. This talk presents an NNTile framework and shows its current results with GPT-based model of different sizes on hardware of a different scale.

Aleksandr Mikhalev
Skoltech
mikhalev@skoltech.ru

## MS28

**Performance Portability of Sparse Matrix Multiple Vector Multiplication on GPUs**

With the recent Department of Energy (DOE) procurement, The National Energy Research Scientific Computing Center (NERSC) Perlmutter, Oak Ridge Leadership Computing Facility (OLCF) Frontier, and Argonne Leadership Computing Facility (ALCF) Aurora supercomputers are accelerated by GPU architectures from NVIDIA, AMD, and Intel, respectively. From that perspective, performance portability across GPU accelerator architectures from various vendors becomes a problem of interest to high-performance computing application developers. These machines are supported by both programming models that are vendor-specific, and portability models that are vendor-agnostic. Therefore, developing performance portable scientific codes for these machines is extremely important to utilize these compute resources efficiently. In this talk, we will demonstrate the challenges to achieve performance portability for large-scale sparse computations. We will particularly focus on sparse matrix vector multiplication (SpMV) and sparse matrix multi-vector multiplication (SpMM) kernels, which constitute a significant portion of sparse computations. We will discuss techniques that can be used to achieve performance portability for SpMV and SpMM kernels on NVIDIA, AMD and Intel based accelerators. To that end, we will explore low-level GPU programming languages such as CUDA and HIP, and high-level GPU portability models such as OpenACC and Kokkos.

Abdullah Alperen
Michigan State University
alperena@msu.edu

## MS28

**Tackling Communication Bottlenecks in Matrix-Free Eigensolvers for Many-Body Localization**

Matrix-free eigensolver for studying the many-body localization (MBL) transition of two-level quantum spin chain models is computationally challenging because the vector space dimension grows exponentially with the physical system size, and averaging over different configurations of the random disorder is needed to obtain relevant statistical behavior. For each eigenvalue problem, eigenvalues from different regions of the spectrum and their corresponding eigenvectors need to be computed. Traditionally, the interior eigenstates for a single eigenvalue problem are computed via the shift-and-invert Lanczos algorithm. Due to the extremely high memory footprint of the LU factorizations, this technique is not well suited for large number of spins. The matrix-free approach does not suffer from these memory bottlenecks. However, its scalability is limited due to communication imbalance and fine-grained irregular transfers. This work presents strategies to reduce this imbalance and optimize communication performance using the consistent space runtime (CSPACER), a runtime designed for application-oriented communication specialization. We will discuss the performance improvement of

the CSPACER-based MBL implementation at scale and compare it to optimized MPI non-blocking two-sided and one-sided RMA implementation variants. The efficiency and effectiveness of the proposed algorithm are demonstrated by computing eigenstates on a massively parallel many-core high-performance computer.

Khaled Z. Ibrahim
Lawrence Berkeley National Laboratory
Berkeley, CA, USA
kzibrahim@lbl.gov

## MS28
### Massively Parallel Relativistic Many-body Method

The fully correlated frequency-independent Dirac-Coulomb-Breit Hamiltonian provides the most accurate description of electron-electron interaction before going to a genuine relativistic quantum electrodynamics theory of many-electron systems. We developed a correlated Dirac-Coulomb-Breit multiconfigurational self-consistent-field method within the frameworks of massively parallel distributed active space. In this approach, the Dirac-Coulomb-Breit Hamiltonian is included variationally in both the mean-field and correlated electron treatment. Benchmark with more than a billion complex-valued spinor determinants is used to demonstrate the parallel efficiency and the computational power of the distributed active space approach. We also analyze the importance of the Breit operator in electron correlation and the rotation between the positive- and negative-orbital space in the no-virtual-pair approximation.

Xiaosong Li
University of Washington
xsli@uw.edu

## MS28
### Modeling Twisted Bi-Layer Graphene by High Performance Computing and Machine Learning

The magic angle has been discovered in Twisted Bilayer graphene (TBG) and electronic structures or quantum material science in the low dimensional layer material has caught the eyes of experts from a number of domains including condensed matter physics, applied analysis, electrical engineering and computational science. The domain particular about twisted bilayer graphene is rapidly being developed. However, there are still difficulties to be solved in the domain, including how to achieve DFT level relaxation and first principle machine learning based modeling for TBG system. We would like to address the modeling problem of twisted bilayer graphene from the perspective of high performance computing and scientific machine learning. The recent progress along these two lines will be presented.

Diyi Liu
University of Minnesota
liu00994@umn.edu

Jalen Harris
Cornell University
jah668@cornell.edu

Lin Lin
University of California, Berkeley
Lawrence Berkeley National Laboratory
linlin@math.berkeley.edu

Chao Yang
Lawrence Berkeley National Lab
cyang@lbl.gov

## MS29
### A Step Beyond Stream-Triggered Communications: AMReXs Stream Object

Networking performance is becoming increasingly critical to full-scale stencil and particle codes on modern supercomputers. Technologies such as stream-triggered and GPU-triggered communications are exciting possibilities to improve performance on accelerator-driven platforms. However, these technologies can be enhanced by exploring holistic solutions that can be extended over the entire application. This investigation introduces an intermediate stream object to minimize the impact of communication synchronization points and present task-based design of applications in a user-friendly package. Built in AMReX, the structured-mesh and particle framework developed at the Lawrence Berkeley National Laboratory, the intermediate stream object launches a unique thread that performs CPU work of corresponding GPU work in an consistent, ordered and optimized manner. This allows the host thread to perform preparatory or meta-data work without synching, while the stream independently completes the critical work path, all in a coding strategy already familiar to GPU programmers: the stream. The design of the intermediate stream object will be explained, including portability to CUDA, HIP and SYCL implementations. Performance results, including overheads and implementations in communication algorithms, will be presented. Finally, broader impact will be discussed with a focus towards next steps and options for universal adoption.

Kevin N. Gott
LBNL
kngott@lbl.gov

Weiqun Zhang
Lawrence Berkeley National Laboratory
Center for Computational Sciences and Engineering
weiqunzhang@lbl.gov

James Dinan
Ohio State University
jdinan@nvidia.com

## MS29
### Performance Anaylsis of E3SM Leveraging Performance Profiling Tools

The physical processes governing the behaviour of the Earth's atmosphere are multi-scale and inter-related. High resolution global models are required to resolve all relevant scales with fidelity. Inevitably this becomes extremely expensive, so the performance of these models limits their predictive power. The Simple Cloud-Resolving

E3SM Atmosphere Model (SCREAM) is a highly optimized GPU-capable code that achieves performance portability by leveraging Kokkos and Ekat. In this presentation I will discuss various metrics of the performance, measured using NVIDIA's profiling tools, in depth. All analysis has been done on NERSC's Perlmutter GPU nodes. Understanding the performance characteristics of the code highlights possible areas of improvement.

Hannah E. Ross
Lawrence Berkeley National Laboratory
hross@lbl.gov

Noel Keen
LBNL
ndkeen@lbl.gov

Peter Caldwell
Lawrence Livermore Nat'l Lab
caldwell19@llnl.gov

## MS29

### Targeted Performance Optimization for Nyx on Hpc Systems

A brief introduction will be given to the different architectures and types of parallel resource constraints relevant to this minisymposium. This talk will focus on performance optimization using the Nyx code as a case study. Nyx is built on top of AMReX, a software framework containing all the functionality to write massively parallel, block-structured adaptive mesh refinement (AMR) applications. The Nyx cosmology code is used to simulate the formation of large scale structure in Ly$\alpha$ simulations of the universe and represents the baryonic matter on the structured mesh and the dark matter as particles. In this talk, we will focus on different performance characteristics of Nyx when constrained by load imbalance and the GPU memory capacity, as well as the various approaches in Nyx to optimizing this across architectures and machines. Of particular interest will be the cumulative effect of interactions such as those between problem characteristics, domain decomposition choice, communication strategies, and compiler choice.

Jean M. Sexton
Lawrence Berkeley National Laboratory
jmsexton@lbl.gov

## MS29

### Modeling Astrophysical Reactive Flows with Castro at the Exascale

Stellar evolution is driven by nuclear reactions. During the normal lives of stars, these reactions can drive convection throughout (parts of) the star. At the end stages of stellar evolution, multiple burning processes can be taking place simultaneously on timescales comparable to (or shorter than) the hydrodynamic timescale. Likewise, stellar remnants can be revived via interactions with their companions, driving explosive burning. Modeling this burning is difficult, and care needs to be taken to ensure that the reactions and hydrodynamics are strongly coupled. The time and lengthscales over which the burning takes place also makes the simulations challenging – high resolution is needed to ensure that the burning is resolved, requiring large computational resources. We'll discuss how we've developed the Castro simulation code to meet these challenges – by building new integration methods to couple reactions and hydrodynamics and leveraging the AMReX adaptive mesh refinement library to make our code performance portable to exascale architectures.

Michael Zingale
Department of Physics and Astronomy
Stony Brook University
michael.zingale@stonybrook.edu

## MS30

### Improving Stability and Performance: Integration of Novel CholeskyQR2 into the ChASE Library

The ChASE library (https://github.com/ChASE-library/ChASE) is used to identify the extremal portion of dense Hermitian eigenvalue problems using the Chebyshev filter. It has been parallelised to support both homogeneous and heterogeneous architectures with distributed memory. In the latest code extension, the original Householder-based QR factorisation used to orthogonalise the filtered vectors has been replaced by the CholeskyQR2 algorithm. This algorithm, which is mainly implemented with level-3 BLAS operations, offers better adaptability to distributed memory systems than the Househoulder QR. However, CholeskyQR2 exhibits numerical instability for ill-conditioned matrices (with a condition number above $10^8$), so one falls back to traditional Householder QR for ill-conditioned matrices. This presentation will focus on ongoing efforts to improve the numerical stability of ChASE QR factorisation. The first part will focus on the parallelisation of the novel CholeskyQR2 algorithm extended with the Gram-Schmidt method for distributed memory and GPU-based systems. The second part focuses on the integration of this novel QR factorisation algorithm into ChASE. The new algorithm shows improved performance and robustness regardless of the matrix condition number. Moreover, due to its design that constructs orthogonalised vectors by panels, it allows for a simpler and more efficient implementation of the updated QR in ChASE, significantly reducing the required flops.

Nenad Mijic
Ruder Bokovic Institute, Croatia
nenad.mijic@irb.hr

Davor Davidovic
Ruder Boskovic Institute
Bijenicka cesta 54, 10000, Zagreb
davor.davidovic@irb.hr

Xinzhe Wu
Jülich Supercomputing Centre
xin.wu@fz-juelich.de

Edoardo A. Di Napoli
Juelich Supercomputing Centre

e.di.napoli@fz-juelich.de

## MS30

### Unite and Learn: An Iterative Approach for AI

An innovative machine learning approach based on Unite and Conquer methods, used in linear algebra, will be presented. In addition to its effectiveness from the point of view of accuracy, the important characteristics of this intrinsically parallel and scalable technique make it well suited to multi-level and heterogeneous parallel and/or distributed architectures. Experimental results, demonstrating the interest of the approach for efficient data analysis in the case of clustering, anomaly detection and road traffic simulation will be presented.

Nahid Emad
LI-PaRAD laboratory and Maison de la Simulation
University of Paris Saclay / Versailles
Nahid.Emad@uvsq.fr

## MS30

### Recent Progress in Complex Moment-based Eigenvalue Solvers

Complex moment-based eigensolvers have been well studied for solving interior eigenvalue problems because of their high parallel efficiency. Typical methods are the FEAST eigensolver and its variants and methods based on Sakurai-Sugiura's approach. This talk will focus on methods based on Sakurai-Sugiura's approach using higher-order complex moments. We introduce recent progress in complex moment-based eigensolvers including a new efficient method and extensions for other related problems.

Akira Imakura
Department of Computer Science
University of Tsukuba
imakura@cs.tsukuba.ac.jp

Tetsuya Sakurai, Yasunori Futamura
University of Tsukuba
sakurai@cs.tsukuba.ac.jp, futamura@cs.tsukuba.ac.jp

## MS30

### Distributed Data Analysis Through Data Collaboration Technique

In this talk, we will present a data collaboration technique for analyzing data distributed across multiple institutions, especially when it contains sensitive information. Instead of sharing the original data, institutions only share intermediate representations, obtained through dimensionality reduction. Importantly, every institution applies its unique transformation function, which remains confidential. We will illustrate some properties of this approach. Some numerical examples show performance of the proposed method.

Tetsuya Sakurai
University of Tsukuba
sakurai@cs.tsukuba.ac.jp

Akira Imakura
Department of Computer Science
University of Tsukuba
imakura@cs.tsukuba.ac.jp

Yasunori Futamura
University of Tsukuba
futamura@cs.tsukuba.ac.jp

## MS31

### Innovative Supercomputing in the Exascale Era by Integration of Simulation/Data/Learning

We propose an innovative method of computational science for sustainable promotion of scientific discovery by supercomputers in the Exascale Era by integrating (Simulation/Data /Learning (S+D+L)), and develop a software platform h3-Open-BDEC for integration of (S+D+L) and evaluate the effects of the integration of on heterogenous supercomputer systems. The h3-Open-BDEC is designed for extracting the maximum performance of the supercomputers with minimum energy consumption. Related activities are described in the talk with future perspectives.

Kengo Nakajima
Information Technology Center, The University of Tokyo
RIKEN Center for Computational Science (R-CCS)
nakajima@cc.u-tokyo.ac.jp

## MS31

### Challenges in Extreme Scale Sparse Linear Algebra Computing

Exascale machines are now available, based on several different arithmetic (from 64-bit to 16-32 bit arithmetics) and using different architectures (with network-on-chip processors and/or with accelerators). Brain-scale applications, from machine learning and AI for example, manipulate huge graphs that lead to very sparse non-symmetric linear algebra problems. Moreover, many supercomputers have been designed primarily for computational science, mainly numerical simulations, not for machine learning and AI. New applications that are maturing after the convergence of big data and HPC to machine learning and AI would probably generate post-exascale computing that will redefine some programming paradigms. End-users and scientists have to face a lot of challenge associated to these evolutions and the increasing size of the data. I review some sparse linear algebra experiments for iterative methods and/or for machine learning. I present some results obtained for linear algebra problems, such as sequence of sparse matrix products and the PageRank method, with respect to the sparsity, and the size of the matrices, on the one hand, and to the number of process and nodes, and the network topologies, on the other hand. Then, I introduce an opensource generators of very large data, allowing to evaluate several methods using very large graph-sparse matrices as data sets for several application evaluations.

Serge G. Petiton
University of Lille, and. CNRS

serge.petiton@univ-lille.fr

## MS31

### Domain Specific Foundation Models - Computational Challenges

The potential role of highly optimized domain specific models is exemplified by prototypes and new software products in pathology and tissue based biomarker development. The required analyses frequently leverage populations of labeled and segmented cells and larger scale microanatomic structures taking into account spatial interrelationships. Domain specific foundation models are being created by us and other groups to support experimentation and development of methods that can be used to steer treatments by predicting response and outcome to alternate treatment regimens. We are developing efficient methods for creating new models using novel self supervised training tasks designed to control attention sparsity patterns along with tasks that involve precise matching tasks that target images such as tissue and satellite imagery. We are also developing methods to efficiently customize existing models via combinations of text and visual prompts and generative AI approaches. Finally, we are developing methods that leverage these models to predict cancer outcome and molecular classification via multi-task training. In this talk, I will describe the methods developed by our group, address computational costs associated with model development, training and evaluation.

Joel Saltz
Department of Biomedical Informatics
Stony Brook University, NY, USA
Stony Brook University, NY, USA

Prateek Prasanna, Tahsin Kurc, Dimitris Samaras
Stony Brook University
prateek.prasanna@stonybrook.edu,
tahsin.kurc@stonybrook.edu, samaras@cs.stonybrook.edu

Saarthak Kapse, Jingwei Zhang
Stony Brook
saarthak.kapse@stonybrook.edu,
jingwezhang@cs.stonybrook.edu¿

Srikar Yellapragada, Jakub Kaczmarzyk
Stony Brook University
srikary@cs.stonybrook.edu,
jakub.kaczmarzyk@stonybrookmedicine.edu

## MS31

### Addressing Extreme-scale Computing and Big Data Demands through Software Ecosystems

Computing systems continue to supply increased performance by leveraging concurrency. GPU, vectorizing CPU, and data flow architecture designs can provide improved performance by executing many operations using large volumes data. New algorithms and software implementations must be utilized to realize performance improvements but these changes are difficult to develop, support, and make portable across the full spectrum of target systems. Software ecosystems–collections of compatible libraries and tools–can address many of the challenges of advanced computing system complexity by providing ready-made, portable functionality that provides optimized implementations of important capabilities that an application code can build upon rather than implement within the application. In this presentation, we discuss ongoing efforts within the US Department of Energy to develop, deliver, and deploy scientific software ecosystems for portable, high-performance execution on a variety of leadership computing systems. These efforts extend to libraries and tools for data science and machine learning with only incremental cost.

Michael A. Heroux
Sandia National Laboratories
St. John's University
maherou@sandia.gov

James Willenbring
Sandia National Laboratories
jmwille@sandia.gov

## MS32

### Reorthogonalized Block Classical Gram-Schmidt Using Two Cholesky-Based TSQR Algorithms

The talk considers the block Gram-Schmidt with reorthogonalization (BCGS2) algorithm discussed in [J. Barlow and A. Smoktunowicz, " Reorthogonalized block classical Gram–Schmidt,' Numerische Mathematik, 123 (2013), pp.395-423] for producing the QR factorization of a matrix $X$ that is partitioned into blocks. A building block operation for BCGS2 is two "tall-skinny QR' factorizations (TSQR) which were assumed to be backward stable orthogonal factorizations such as Householder or Givens QR. However, that assumption of backward stability excludes some possible TSQR algorithms, in particular, the first of these two factorizations could be done using a mixed precision Cholesky TSQR algorithm from [I. Yamazaki, S. Tomov, and J. Dongarra, "Mixed-precision Cholesky QR factorization and its case studies on multicore CPU with multiple GPUs", SIAM J. Sci. Computing, 37 (2015), pp.C307-C330]. The second TSQR could be done with a working precision Cholesky-based QR decomposition. It is shown that these significantly weaker stability conditions such as those satisfied by these two TSQR algorithms are sufficient for BCGS2 to produce a conditionally backward stable factorization.

Jesse L. Barlow
Penn State University
Dept of Computer Science & Eng
barlow@cse.psu.edu

## MS32

### QRCP of a Tall-skinny Matrix by a Cholesky QR Type Algorithm

QR factorization with column pivoting (QRCP) is a fundamental matrix factorization and has important applications including low-rank approximation and rank determination. Recently, it has been reported that the QR factorization (without pivoting) of a tall and skinny matrix is efficiently

computed by Cholesky QR type algorithms. Motivated by this, we aim to develop a Cholesky QR type algorithm that computes the QRCP of a tall-skinny matrix. In the algorithm development, the main difficulty lies in the accurate selection of columns (pivots) under the effect of rounding error, and our solution is to determine reliable pivots in an iterative manner. Our algorithm has the same structure as existing Cholesky type algorithms, which is suitable for recent computer architectures; almost all computations can be done in Level-3 BLAS routines such as GEMM, and the number of required collective communication in distributed parallel computing is $O(1)$, i.e., it can be regarded as communication avoiding. In numerical experiments, our algorithm provides QRCP as accurate as those obtained by conventional algorithms. Furthermore, our algorithm is faster than conventional algorithms in both single node and distributed parallel environments.

Takeshi Fukaya
Hokkaido University
fukaya@iic.hokudai.ac.jp

Yuji Nakatsukasa
University of Oxford
nakatsukasa@maths.ox.ac.uk

Yusaku Yamamoto
The University of Electro-Communications, Japan
yusaku.yamamoto@uec.ac.jp

## MS32
### Two-step Block Gram-Schmidt in s-step GMRES

Compared to the standard GMRES, its s-step variant provides the potential to reduce the communication cost of orthogonalizing a set of s Krylov vectors. However, these Krylov vectors generated by the matrix-powers kernel become increasingly ill-conditioned as a larger value of s is used. As a result, to maintain stability in practice, a conservatively small step size is used, which can limit the performance gain that s-step GMRES can provide. In this talk, we explore a "two-step' Gram-Schmidt, where the set of s Krylov vectors are first "pre-processed' to maintain stability, while their orthogonalization is delayed until enough vectors are generated to obtain better performance. To maintain the conditioning of the Krylov vectors, we look at Block Classical Gram-Schmidt (without reorthogonalization) and random-sketching techniques.

Ichitaro Yamazaki
Sandia National Laboratories
iyamaza@sandia.gov

Andrew J. Higgins
Temple University, U.S.
andrew.higgins@temple.edu

Andrew J. Higgins
Temple University
andrew.higgins@temple.edu

Erik G. Boman
Center for Computing Research
Sandia National Labs
egboman@sandia.gov

Daniel B. B. Szyld
Temple University
Department of Mathematics
szyld@temple.edu

## MS33
### Streamlining Massively-Parallel Work-Flows with the Fierro Mechanics Code

We present recent research for efficient end-to-end work flows with the open-source massively parallel FIERRO mechanics code. FIERRO is a C++ code to aid: (a) modeling efforts, (b) developing novel numerical methods, (c) investigating performance portability, and (d) testing of user-defined models related to quasi-static and dynamic problems involving fluids, solids, and gasses. FIERRO is portable across diverse types of high performance computing (HPC) machines. Fine-grain parallelism is used for on-node parallelism with multi-core CPUs and GPUs, while the message passing interface (MPI) is used for distributed, massive parallelism. FIERRO contains implicit Lagrangian finite element methods to simulate quasi-static problems, such as calculating the stress in a metal beam, the steady-state heat transfer, or the thermal-mechanical response of a part. Novel density-based multi-physics topology optimization approaches are offered in the code. For material dynamics applications – including simulating fluid, gas, and solid dynamics – FIERRO contains a diverse suite of Lagrangian methods including arbitrary-order continuous and discontinuous finite element methods. Given the extensive capabilities of the FIERRO code, and a growing user-base, it is imperative to develop simple, straightforward user interfaces. This talk discusses the research on workflow tools to efficiently use FIERRO and analyze large-scale simulation results.

Sarah Hankins, D Dunning, Kevin Welsh, Nathaniel Morgan
Los Alamos National Laboratory
shankins@lanl.gov,                daniel.dunning@lanl.gov,
kevin.welsh@lanl.gov, nmorgan@lanl.gov

## MS33
### Simulation Workflow Capabilities of the Moose Framework

One of the most overlooked parts of developing successful scientific software is the day-to-day workflow of users and developers. The MOOSE (Multiphysics Object Oriented Simulation Environment) project has always sought to balance adding capabilities and enabling use of the software. Over 15 years of development have gone into creating streamlined workflows for both application developers and end-users. For developers, a bespoke build-system simplifies the compilation of complicated multiphysics systems. Custom testing and continuous integration systems run over 50 million tests per week. A Conda-based binary distribution system enables developers to instantly create a working development environment and keep it up to date as hundreds of changes are made to the ecosystem each week. Users of MOOSE-based codes benefit from multiple work-

flow enhancing capabilities. All MOOSE-based codes use a similar input syntax and automatically generate documentation. Built-in mesh generation, postprocessing, and statistical analysis greatly speed up AI/ML and uncertainty quantification. Multiple graphical user interfaces provide a streamlined user experience. This talk will detail these capabilities and their impact on the growth and development of the overall ecosystem.

Derek R. Gaston
Idaho National Laboratory
derek.gaston@inl.gov

Cody J. Permann
Advanced Modeling and Simulation
Idaho National Laboratory
cody.permann@inl.gov

Alexander Lindsay, Logan Harbour, Jason Miller, Casey Icenhour
Idaho National Laboratory
alexander.lindsay@inl.gov, logan.harbour@inl.gov, jason.miller@inl.gov, casey.icenhour@inl.gov

## MS33
### Instrumenting Simulations with Computational Model Builder

Large-scale simulations are seldom the result of a single person working alone, but rather a team of experts in different areas. As these simulations are prepared, people with different skills need a collaborative tool to shepherd the simulation through the setup, simulation, and post-processing workflow. Computational Model Builder (CMB) is such an end-to-end workflow application built on the Simulation Modeling Tool Kit (SMTK) framework. In this presentation, we discuss the scheme it adopts for describing simulation input parameters which are then exported to input decks for runs on large parallel machines. By providing a language for describing the simulation, SMTK can validate parameters that might otherwise lead to wasted time. Similarly, our tool lets users graphically associate materials and boundary conditions to the simulation's geometric domain(s) to prevent costly errors. CMB's user interface is built on ParaView, so simulations can also be inspected remotely while they are running – from the same application used to prepare their input. In addition to describing the capabilities of CMB, we illustrate its use on applications such as linear accelerator design, anatomical modeling, hydrological modeling, nuclear reactor design, and wave-energy capture.

Robert O'Bara
Kitware
bob.obara@kitware.com

John Tourtellott, David Thompson
Kitware, Inc.
john.tourtelloutt@kitware.com,
david.thompson@kitware.com

## MS33
### Simput and Trame for Web-Based Simulation

### Workflows

Large parallel simulations can be difficult to prepare; improperly configuring a simulation can be costly and scalable simulations often have complex input grammars that make failure easy. Often, a particular application will use only a small fraction of a simulation's input parameters. Trame enables rapid development of user interfaces (UIs) to address these situations with minimal effort. Trame relies on scalable, high-level components quickly assembled into a UI, handling complex interactions – such as between UI elements, charts, 3D visualizations, and more – seamlessly with minimal friction by leveraging best-of-class Python libraries like matplotlib, plotly, VTK, and ParaView. The framework organically provides a ubiquitous deployment approach that enables desktop, client/server, cloud, HPC, and Jupyter scenarios from a single code base. For simulation workflows that require numerous parameters with complex interdependencies and constraints, Trame integrates Simput, which handles form generation and data persistence from lightweight model definitions. Leveraging Trame and Simput, one can focus on the data rather than the presentation layer. We demonstrate this approach by describing the development of the ArrowFlow microworkflow product based on Trame, Simput, and M-Star's simulation; by targeting a specific simulation area (chemical reactors for pharmaceuticals), the UI was simple and inexpensive to develop but easy to customize as needed.

Sebastien Jourdain
Kitware, Inc.
sebastien.jourdain@kitware.com

Patrick O'Leary
Kitware, Inc., USA
patrick.oleary@kitware.com

David Thompson
Kitware, Inc.
david.thompson@kitware.com

## MS34
### Nonlinear Multifidelity Optimization for Inverse Electrophysiology

We present and discuss a multi-fidelity classifier for the evaluation of the atrial fibrillation inducibility in patient specific models. Our classifier classifier is based on a multi-fidelity approach, which combinse accurate high fidelity models with fast low fidelity models, thereby efficiently reducing the computational cost of the personalized classification process. Here, we combine a large-scale high fidelity monodomain model and a low fidelity model based on a coarser discretization. Moreover, we present a low fidelity eikonal model to further improve the computational efficiency in simulations of atrial fibrillation. For this purpose, the standard eikonal algorithm has to be adapted to account for the re-excitability of the tissue. Numerical experiments show that the multi-fidelity classifier is more efficient than the single-fidelity classifiers and that the eikonal low fidelity model provides promising results in the fast approximation of the monodomain high fidelity

model.

Rolf Krause
Università della Svizzera italiana, Switzerland
rolf.krause@usi.ch

## MS34

### Reconstructing Cardiac Electrical Excitations from Optical Mapping Recordings

The reconstruction of electrical excitation patterns through the unobserved depth of the tissue is essential to realizing the potential of computational models in cardiac medicine. We have utilized experimental optical-mapping recordings of cardiac electrical excitation on the epicardial and endocardial surfaces of a canine ventricle as observations directing a local ensemble transform Kalman Filter (LETKF) data assimilation scheme. We demonstrate that the inclusion of explicit information about the stimulation protocol can marginally improve the confidence of the ensemble reconstruction and the reliability of the assimilation over time. Approximation error is addressed at both the observation and modeling stages, through the uncertainty of observations and the specification of the model used in the assimilation ensemble. We find that incorporating additional information from the observations into the dynamical model itself (in the case of stimulus and stochastic currents) has a marginal improvement on the reconstruction accuracy over a fully autonomous model, while complicating the model itself and thus introducing potential for new types of model error. That the inclusion of explicit modeling information has negligible to negative effects on the reconstruction implies the need for new avenues for optimization of data assimilation schemes applied to cardiac electrical excitation.

Christopher Marcotte
Computer Science
Durham University
christopher.marcotte@durham.ac.uk

Matthew J. Hoffman
Rochester Institute of Technology
mjhsma@rit.edu

Elizabeth M. Cherry
Georgia Tech
School of Computational Science and Engineering
elizabeth.cherry@gatech.edu

Flavio H. Fenton
Georgia Institute of Technology
flavio.fenton@physics.gatech.edu

## MS34

### Automated Computation of Strength-Interval Curves Using Actors

Strength-interval (SI) curves are used to quantify the relationship between the response of excitable tissue as a function of the strength and duration of an electrical stimulus. In the context of cardiac electrophysiology, SI curves characterize the refractoriness of cardiac tissue as a function of inter-stimulus interval length. Although usually collected experimentally, this type of information can now also be obtained computationally. However, the computational generation of SI curves can be labor intensive and time consuming due to its iterative nature, the number and size of computations required, and the amount of researcher intervention involved. In this presentation, we demonstrate how the actor model of concurrent programming can be used to automate the process of SI curve generation, relieving much of the burden on the researcher while maximizing use of the available computational resources. Computational resource utilization is optimized by the dynamic monitoring and assessment of the overall benefit of each actor's momentary resource usage. The automatically generated SI curves are produced in significantly less time than manually generated ones with the number of same data points. Additionally, the actor-determined threshold stimuli were more precise compared to those that were manually generated. A newly proposed concurrent actor-based look-ahead bisection method for event location is also described.

Raymond J. Spiteri
University of Saskatchewan, Canada
Department of Computer Science
spiteri@cs.usask.ca

Kyle Klenk, Joyce Reimer
Department of Computer Science
University of Saskatchewan
kyle.klenk@usask.ca, joyce.reimer@usask.ca

## MS34

### Enabling Cell-Based Simulations of Cardiac Electrophysiology with High-Performance Computing.

Cell-based simulations represent a new modeling approach to studying cardiac electrophysiology. In this approach, the individual cardiac cells and their sub-cellular details are resolved, distinguishing between the intra-cellular and extra-cellular spaces, while also detailing the membrane space in between. A new challenge arises with the resulting extreme amounts of computation. We will use several techniques of high-performance computing in this context, with the purpose of enabling efficient cell-based simulations on supercomputers. Specifically, we will investigate a new problem of "3-space" mesh partitioning that has the overall objective of minimizing the resulting communication overhead, while ensuring a satisfactory level of load balancing between the processors when numerically solving all the three subproblems. The opportunities of overlapping computation with communication will also be investigated. Numerical experiments will be presented to analyze the impact of different schemes for solving the new mesh-partitioning problem, as well as the effect of communication-computation overlap.

Xing Cai
Simula Research Laboratory
1325 Lysaker, Norway
xingca@simula.no

James D. Trotter
Simula Research Laboratory

james@simula.no

## MS35

### Eigensolvers for Configuration Interaction Calculations of Atomic Nuclei on Exascale Systems

Ab initio calculations of the structure of atomic nuclei using the Configuration Interaction (CI) approach require computing the lowest eigenvalues and eigenvectors of a very large but extremely sparse symmetric matrix. The size of these matrices can be in the (tens of) billions, which require efficient algorithms and implementations on exascale systems. The Lanczos algorithm with orthonormalization is generally the most robust method, but by employing a suitable set of initial starting vectors for the Locally Optimal Block Preconditioned Conjugate Gradient (LOBPCG) solver or for the Residual Minimization Method accelerated by Direct Inversion of Iterative Subspace (RMM-DIIS) one can significantly reduce the number of iterations necessary for a converged set of eigenvectors. Furthermore, an advantage of the RMM-DIIS algorithm is that one can selectively refine specific eigenvectors of interest. These methods have been implemented in Many-Fermion Dynamicsnuclear (MFDn), which is a hybrid parallel MPI + OpenMP Fortran code for ab-initio nuclear structure calculations, with OpenACC or OpenMP for target offloading. For large-scale calculations, data movement and in particular inter-node communication becomes a bottleneck; we therefore also have a matrix-free implementation of the matrix-vector multiplication. This matrix-free implementation is most beneficial for block algorithms, and in principle allows for calculations that are too large to fit in memory.

Pieter Maris
Iowa State University
pmaris@iastate.edu

## MS35

### Dynamical Downfolding and Formulation of Effective Hamiltonians

Quantum many-body calculations, especially those addressing the time evolution of electrons in realistic systems, necessitate an effective representation that reduces computational resources. In practice, the original problem, which is intractable, is divided into an explicitly treated subspace and an approximately described environment. The two are typically separated based on the energy of individual excited states and their nature (e.g., originating in localized states), and only the subspace solution is sought. This reduced representation corresponds to nonlinear eigenvalue problems with space-time non-local operators (capturing the renormalization of one- and two-body interactions). I will discuss the techniques for studying such excited states and outline the challenges associated with these approaches. Further, I will comment on the possible limitations and the role of connection between the states (eigenvectors) obtained for the subsystem and the original many-body problem.

Vojtech Vlcek
University of California Santa Barbara
vlcek@ucsb.edu

## MS35

### Massively Parallel Selected Configuration Interaction for First-Principles Quantum Chemistry Applications

The many-body simulation of quantum systems is an active field of research that involves several different methods targeting various computing platforms. Many methods commonly employed, particularly coupled cluster methods, have been adapted to leverage the latest advances in modern high-performance computing. Selected configuration interaction (sCI) methods have seen extensive usage and development in recent years. However, the development of sCI methods targeting massively parallel resources has been explored only in a few research works. Here, we present a parallel, distributed memory implementation of the adaptive sampling configuration interaction approach (ASCI) for sCI. In particular, we will address the key concerns pertaining to the parallelization of the determinant search and selection, Hamiltonian formation, and the variational eigenvalue calculation for the ASCI method

David B. Williams-Young
Lawrence Berkeley National Laboratory
dbwy@lbl.gov

Norman Tubman
NASA Ames
norm.m.tubman@gmail.com

Bert de Jong
Lawrence Berkeley National Lab
wadejong@lbl.gov

## MS36

### Code Generation for Matrix-Free Finite Element Methods on Hybrid Tetrahedral Grids

This talk presents a code generation approach to flexibly generate optimized kernels for matrix-free finite element methods on block-structured tetrahedral grids. Our approach enables automated optimizations that heavily exploit the underlying grid structure and are agnostic to the targeted bilinear form. We give an overview of the pipeline and analyze various optimization techniques through the application of resource-based performance models. Finally, we demonstrate the extreme scalability of the generated matrix-free methods via applications to linear systems with trillions ($10^{12}$) of unknowns.

Nils Kohl
University of Erlangen-Nuremberg (FAU)
nils.kohl@fau.de

Fabian Böhm
Friedrich-Alexander-Universität Erlangen-Nürnberg
Fabian.boehm@fau.de

Daniel Bauer
University of Erlangen-Nuremberg (FAU)
daniel.j.bauer@fau.de

Ulrich Rüde
Friedrich-Alexander-Universität Erlangen-Nürnberg
(FAU)
ulrich.ruede@fau.de

## MS36

### Implementing a GPU-Resident BDDC Preconditioner for Cardiac Simulations with Sub-Cellular Resolution

While current simulations of the electrophysiology of the human heart mostly consider hundreds of cells per model element in the monodomain and bidomain model, these models are not powerful enough to take into account the heart's individual cells. This is necessary in order to understand effects structural muscle damage that is caused by aging or heart disease has on the heartbeat. In the European MICROCARD project we work on such simulations with the Extracellular, Membrane and Intracellular (EMI) model. The model's high resolution and the resulting linear systems pose a computational challenge that requires tailored preconditioners and solvers capable of effectively leveraging modern exascale hardware. In this talk, we present our efforts on implementing a GPU resident BDDC preconditioner in the open source library Ginkgo that we aim to leverage in these simulations.

Fritz Göbel, Marcel Koch, Terry Cojean
Karlsruhe Institute of Technology
fritz.goebel@kit.edu, marcel.koch@kit.edu,
terry.cojean@kit.edu

Hartwig Anzt
University of Tennessee, U.S.
hanzt@icl.utk.edu

## MS36

### Scalable Algebraic Multigrid Methods for High-fidelity Flow Problems

The Algebraic Multigrid (AMG) method has over the years developed into an efficient tool for solving unstructured linear systems. Solving large unstructured problems has been a key motivation for devising a scalable and efficient AMG method. Despite some success, the key part of the AMG algorithm; the coarsening step, is far from trivial to parallelize efficiently. We here introduce a novel parallelization of the inherently sequential Ruge-Stben coarsening algorithm that retains most of the original method's good interpolation and convergence properties. Our parallelization is based on the Partitioned Global Address Space (PGAS) abstraction, which significantly simplifies the parallelization compared to traditional message passing based implementations. The new coarsening algorithm has been integrated into the high-fidelity flow solver Neko, using a hybrid p-MG/AMG method for efficiently solving the pressure equation in direct numerical simulation of turbulent flow.

Niclas Jansson
KTH Royal Institute of Technology
njansson@kth.se

## MS37

### Systematic Application-Derived Communication System Performance Analysis and Benchmarking

Proxy applications, mini-applications, and micro-benchmarks have become the de-facto standard for systematically studying the performance of high-end computing systems. Unfortunately, few of these codes have algorithms or input decks that are representative of the complexity of communication present in real codes. While a wide range of mini-applications and proxy applications seek to reproduce certain application characteristics, input problems are often simplified in ways that trivialize the resulting communication patterns. The situation is even more dire for understanding the communication behavior of coupled or multi-scale codes; we are not aware of any benchmark that attempts to reproduce any aspect of the communication behavior of such codes. In this talk, I describe the approach we are taking to address this issue. First, we are developing measurement, characterization, and benchmarking techniques that can extract and replay non-trivial communication patterns from production HPC applications on production workloads without the complexity of tracing; this allows us to easily study the impact of communication system changes on these workloads. Second, we are identifying key communication patterns and application spaces not covered by current benchmarks and developing new mini-applications to drive research in this direction. To illustrate our advances in this space, I will present the results of communication patterns and statistics extracted from the LANL xRage production code, and describe the design and implementation of Beatnik, a high-performance parallel interface mini-application that captures key elements of codes with complex particle/mesh data couplings and redistribution and is a first step toward developing true coupled-code mini-applications.

Patrick G. Bridges
University of New Mexico
bridges@cs.unm.edu

## MS37

### Partitioned Communication and the Future of Application Design

Scientific applications continue to be run on increasingly complex systems. Between an increased level of CPU complexity and the use of accelerators such as GPUs; multithreaded and accelerator-driven communication is becoming increasingly necessary for both performance and usability. However, the current MPI Send Recv model is not optimized for these use cases. In this talk, I will present some the exploratory work Sandia has done on benchmarking future application paradigms. This work, MiniMOD and CMB, ranges from design and implementation to novel performance evaluation and optimization. MiniMOD is a modular communication framework, designed for benchmarking purposes. It allows for direct comparisons to be made between different fine-grained communication implementations. The CMB (Configurable Micro Benchmark)

is a benchmark designed to allow for the exploration of different application profiles including message size, computational time, and thread completion distribution. The combination of these two allow us to explore not only what future application might look like, but also how they will perform. In this talk I will present how these and other advances allow us to explore the future of application design.

Matthew Dosanj
Sandia National Laboratories
mdosanj@sandia.gov

## MS37

### Developing "Useful" Proxy Applications - Problems, Solutions, and Problems

Modern scientific applications are complicated and require coordination of several components. Proxy application driven software-hardware co-design plays a vital role in driving innovation among the developments of applications, software infrastructure and hardware architecture. Proxy applications are self-contained and simplified codes that are intended to model the performance-critical computations within applications. While there is disagreement in the HPC community on the mechanisms of construction of the proxy applications, there is a strong consensus on their positive impact in co-design. The need for achieving sustainable strong scalability on future extreme-scale heterogeneous systems drives the considerations for designing appropriate proxy applications and abstractions. We need a multi-pronged strategy to 1) develop/utilize high-level distributed-memory programming abstractions for rapid prototyping of application scenarios, 2) devise efficient heuristics and methods to optimize data movement for application scenarios (especially sparse and irregular applications that are memory-access driven), and 3) quantitative evaluation to assess systems performance. In this talk, we will discuss these aspects of proxy applications driven co-design.

Sayan Ghosh
Pacific Northwest National Laboratory
sayan.ghosh@pnnl.gov

## MS37

### Livermores Perspective on Functional Reproducibility, Automation, and Community Collaboration in HPC Benchmarking

We use benchmarking to communicate our computational workloads, and verify performance of the delivered HPC systems. Yet every step in HPC benchmarking is manual: even if a port of the benchmark to a given hardware exists, building and running is different on every system. The high barrier to entry is hampering reproducibility of the benchmarks and interaction between HPC centers, vendors, and researchers. We propose collaborative continuous benchmarking to enable functional reproducibility, automation, and community collaboration in HPC benchmarking. This talk will describe our initial implementation of Benchpark, and open the discussion on how collaborative benchmarking will help overcome the human bottleneck in HPC benchmarking, enabling better evaluation of our systems and a more productive collaboration within the HPC community.

Olga Pearce
Lawrence Livermore National Lab
pearce8@llnl.gov

## MS38

### Multifrontal Solver with Low-Rank Compression on CPUs and GPUs

This talk presents an approximate sparse multifrontal solver, used as a preconditioner for GMRES. We leverage hierarchical and non-hierarchical matrices, using low-rank compression methods, to reduce the asymptotic complexity of computation and memory usage. The multifrontal solver can combine multiple rank structured formats, such as block low rank and hierarchically off-diagonal butterfly. In particular, we apply the hierarchically off-diagonal butterfly compression to large sized fronts of the multifrontal method and apply the block low rank compression to medium sized fronts. This combination reduces both the asymptotic complexity as well as the runtime for medium scale problem sizes. We illustrate the performance of this new preconditioner on CPUs as well as modern GPU architectures. We consider a range of industrial and academic applications: high frequency Helmholtz and Maxwell, incompressible Navier-Stokes, etc.

Lisa Claus
Lawrence Berkeley National Laboratory
lclaus@lbl.gov

Pieter Ghysels
Lawrence Berkeley National Laboratory
Computational Research Division
pghysels@lbl.gov

Yang Liu, Xiaoye (Sherry) Li
Lawrence Berkeley National Laboratory
liuyangzhuan@lbl.gov, xsli@lbl.gov

## MS38

### Fast Direct Solution of EM Scattering Problems with Hierarchical Matrices

In this work, H-Matrix acceleration of the Locally Corrected Nystrom (LCN) solution is represented. Generally, H-Matrix is used to approximate a dense matrix using a data-sparse representation. To do this, the matrix should be partitioned into blocks and for those blocks satisfying the admissibility criteria, low-rank decomposition is implemented. In this new computational framework, H-LU decomposition and H-back-substitution are employed to accelerate direct solution of the pertinent. Also, Adaptive Cross Approximation (ACA) and error-controlled H-matrix arithmetic provides error-controlled fast direct solution to the matrix equations by compressing the H-matrix blocks. Using the proposed approach, it is demonstrated that smooth objects of electrically moderate size (up to approximately 100 wavelengths) can be solved with O(NlogN) CPU time and memory consumption. The key factor of H-

Matrix is in subdividing of the matrix using the recursive blocks. It means, the area needs to be partitioned into similar sub-areas and it is based on geometry of the set of row and column indices of the matrix. This is a vital step in physical distances and interaction between degrees of freedom. The conference will present numerical results demonstrating $O(h^p)$ error behavior as well as CPU and memory complexity scaling for the H-Matrix acceleration of the LCN method.

Omid Babazadeh
University of Manitoba
babazado@myumanitoba.ca

Jin Hu
University of Southern California
jinhu@usc.edu

Emrah Sever
emrahsever@aselsan.com.tr
emrahsever@gtu.edu.tr

Ian Jeffrey
University of Manitoba
ian.jeffrey@umanitoba.ca

Constantine Sideris
University of Southern California
csideris@usc.edu

Vladimir Okhmatovski
University of Manitoba
vladimir.okhmatovski@umanitoba.ca

## MS38

### A Mixed Sparse-Dense BLR Solver for Electromagnetics

Element-by-element preconditioners were an active area of research in the 80s and 90s, and they found some success for problems arising from Finite Element discretizations, in particular in structural mechanics and fluid dynamics (e.g., the "EBE" preconditioner of Hughes, Levit, and Winget). Here we consider problems arising from Boundary Element Methods, in particular the discretization of Maxwell's equations in electromagnetism. We focus on an iterative solver using a matrix splitting $A = M + N$, where $M$ is used as a preconditioner. Each iteration requires a solve with $M$ and a matrix-vector multiply with $N$. $M$ is built from the elemental matrices associated with pairs of adjacent elements. The degree of adjacency (e.g., sharing a vertex, sharing an edge, etc.) gives different preconditioners. $M$ is sparse and can be factored with a direct method, possibly accelerated with low-rank techniques. $N$ holds the remainder of the elemental matrices; it is dense and can be compressed to accelerate the matrix-vector products. We use the Block Low-Rank approach (BLR). In the BLR approach, a given dense matrix (or submatrix, in the sparse case) is partitioned into blocks following a simple, flat tiling; off-diagonal blocks are compressed into low-rank form using a rank-revealing factorization, which reduces storage and the cost of operating with the matrix. We demonstrate results for large industrial problems coming from the LS-DYNA multiphysics software.

Francois-Henry Rouet
Ansys, Inc
francois-henry.rouet@ansys.com

Cleve Ashcraft, Pierre L'Eplattenier
Ansys, Inc.
cleve.ashcraft@ansys.com, pierre.leplattenier@ansys.com

## MS38

### A Stable Matrix Version of Some Fast Transforms Based on Sum-of-Exponentials Expansions

We give a stable hierarchical algorithm of some fast transforms for evaluating the matrix-vector product to some discretized kernel matrices, such as the Gaussian kernel the and electrostatic potential kernel. Existing work shows that, with the help of sum-of-exponentials expansions, these kernel matrices can be quickly approximated into some structured matrices. However, it turns out that operations with such structured forms are performed sequentially and have stability risks. We analyze relevant stability and use the sum-of-exponentials expansions to construct a hierarchical transform with guaranteed stability and much better scalability. The error propagation is shown to be proportional to a low-degree power of $\log(n)$, where $n$ is the matrix size.

Chenyang Cao, Jianlin Xia
Purdue University
cao302@purdue.edu, xiaj@purdue.edu

## MS39

### A Quantum-walk Inspired Ansatz for Variational Quantum Algorithm

Continuous-time quantum walks (CTQWs) on dynamic graphs are a universal model of computation that have striking parallels to the quantum approximate optimization algorithm (QAOA), a variational quantum algorithm typically used to solve combinatorial optimization problems. In this talk, we first introduce a dictionary that efficiently converts CTQWs on dynamic graphs to quantum circuit elements. We then use this dictionary to show how quantum walks on dynamic graphs can be used to develop a QAOA ansatz for optimization problems. Finally, we show how this ansatz can be used to prepare particular quantum states.

Rebekah Herrman
University of Tennessee Knoxville
rherrma2@utk.edu

Alvin Gonzales
Argonne National Laboratory
agonza@siu.edu

Colin Campbell, Teague Tomesh
Infleqtion
colin.campbell@infleqtion.com,
teague.tomesh@infleqtion.com

Ji Liu, Zain Saleem
Argonne National Laboratory
ji.liu@anl.gov, zsaleem@anl.gov

## MS39

### Constructing Optimal Contraction Trees for Tensor Network Quantum Circuit Simulation

One of the key problems in tensor network based quantum circuit simulation is the construction of a contraction tree which minimizes the cost of the simulation, where the cost can be expressed in the number of operations as a proxy for the simulation running time. This same problem arises in a variety of application areas, such as combinatorial scientific computing, marginalization in probabilistic graphical models, and solving constraint satisfaction problems. In this paper, we reduce the computationally hard portion of this problem to one of graph linear ordering, and demonstrate how existing approaches in this area can be utilized to achieve results up to several orders of magnitude better than existing state of the art methods for the same running time. To do so, we introduce a novel polynomial time algorithm for constructing an optimal contraction tree from a given order. Furthermore, we introduce a fast and high quality linear ordering solver, and demonstrate its applicability as a heuristic for providing orderings for contraction trees. Finally, we compare our solver with competing methods for constructing contraction trees in quantum circuit simulation on a collection of randomly generated Quantum Approximate Optimization Algorithm Max Cut circuits and show that our method achieves superior results on a majority of tested quantum circuits.

Cameron Ibrahim
University of Delaware
cibrahim@udel.edu

Danylo Lykov
Argonne National Laboratory
dlykov@anl.gov

Zichang He
University of California, Santa Barbara
zichanghe@ucsb.edu

Yuri Alexeev
Argonne National Laboratory
yuri@alcf.anl.gov

Ilya Safro
University of Delaware
isafro@udel.edu

## MS39

### Novel Inexact Feasible Quantum Interior Point Methods for Linear Optimization

The use of quantum computing to accelerate complex optimization problems is a fast growwing research field. This paper applies Quantum Linear System Algorithms (QLSAs) to Newton systems within Interior Point Methods (IPMs) to take advantage of quantum speedup in solving Linear Optimization (LO) problems. Due to their inexact nature, QLSAs can be applied only to inexact variants of IPMs. Existing IPMs with inexact Newton directions are infeasible methods due to the inexact nature of their computations. We proposes an Inexact-Feasible IPM (IF-IPM) for LO problems, using a novel linear system to generate inexact but feasible steps. We show that this method has $(\sqrt{n}L)$ iteration complexity, analogous to the best exact IPMs, where $n$ is the number of variables and $L$ is the binary length of the input data. Moreover, we examine how QLSAs can efficiently solve the proposed system in an iterative refinement (IR) scheme to find the exact solution without excessive calls to QLSAs. We show that the proposed IR-IF-IPM can also be helpful to mitigate the impact of the condition number when a classical iterative method, such as a Conjugate Gradient method or a quantum solver is used at iterations of IPMs. After applying the proposed IF-IPM to the self-dual embedding formulation, we investigate the proposed IF-IPM's efficiency using the QISKIT simulator of QLSA.

Tamas Terlaky
Lehigh University
Department of industrial and Systems Engineering
terlaky@lehigh.edu

Mohammadhossein Mohammadisiahroudi
Industrial & System Engineering Dep, Lehigh University
625 Montclair Ave, Bethlehem, PA, USA, 18015
mom219@lehigh.edu

Ramin Fakhimi
Industrial & System Engineering Dep, Lehigh University
raf318@lehigh.edu

Zeguan Wu
Lehigh University, PA, USA
zew220@lehigh.edu

## MS39

### Hybrid Quantum-Classical Multilevel Algorithm for the Max-Cut Problem

Combinatorial optimization is one of the fields where near term quantum devices are being utilized with hybrid quantum-classical algorithms to demonstrate potentially practical applications of quantum computing. One of the most well studied problems in combinatorial optimization is the Max-Cut problem. The problem is also highly relevant to quantum and other types of post Moore architectures due to its similarity with the Ising model and other reasons. We introduce a scalable hybrid multilevel approach to solve large instances of Max-Cut using both classical only solvers and quantum approximate optimization algorithm (QAOA). We compare the results of our solver to existing state of the art large-scale Max-Cut solvers. We demonstrate excellent performance of both classical and hybrid quantum-classical approaches and show that using QAOA within our framework is comparable to classical approaches.

Anthony Angone
University of Delaware
aangone@udel.edu

Xiaoyuan Liu
Fujitsu Research USA
xliu@fujitsu.com

Ruslan Shaydulin
JPMorgan Chase
ruslan.shaydulin@jpmchase.com

Ilya Safro
University of Delaware
isafro@udel.edu

**MS40**

**Design Patterns and Performance Analysis of Polymorphism in Multiphysics FE Assembly on GPU**

The computation of elemental system matrices and right-hand-side vectors and their assembly into sparse linear algebra data structures is a key component of FE codes. This computation naturally exposes parallelism (loops over regions of PDEs with different coefficients or of different type, loops over elements, loops over integration points, etc). This parallelism can be exploited by implementing the computation in terms of functors that are executed in *Kokkos* parallel regions. In this talk, we focus on such a *Kokkos* implementation in the context of a multiphysics FE code. Specifically, we investigate polymorphic design patterns to isolate commonality across physical applications and provide interfaces to accommodate physical application specificities (differences in data structures used for scalar-valued versus vector-valued FE basis functions, differences in actions of differential operators of different types, etc). The use of polymorphism on device can have performance impacts and is challenging in terms of code generation. We will discuss how we used and compared static and dynamic polymorphic design patterns, and analyzed the performance impact, especially on GPU (virtual function call overhead, compiler optimisations, etc). The design patterns will be illustrated and the performance will be evaluated in the context of a computational electromagnetism simulation relevant to diffraction gratings.

Maarten Arnst, Romin Tomasetti
Universite de Liege
maarten.arnst@uliege.be, romin.tomasetti@uliege.be

**MS40**

**Performance Portable Automatic Differentiation Tools for Finite Element Assembly on Next-generation Architectures**

Supporting scalable and performant finite element simulations across next generation heterogeneous architectures can add significant code complexity. This presentation will discuss the design of automatic differentiation tools for finite element assembly. Performance portability is achieved via the Kokkos programming model. Automatic differentiation is implemented using operator overloading with expression templates. The assembly routines are templated on the scalar type. Performance results will be shown for NVIDIA GPU and Intel Haswell architectures. Examples will be shown for magnetohydrodynamics and multi-fluid

plasma simulations.

Roger Pawlowski
Multiphysics Simulation Technologies Dept.
Sandia National Laboratories
rppawlo@sandia.gov

Eric Phipps
Sandia National Laboratories
Center for Computing Research
etphipp@sandia.gov

Christian Trott
Sandia National Laboratories
crtrott@sandia.gov

John N. N. Shadid
Sandia National Laboratories
Albuquerque, NM
jnshadi@sandia.gov

Eric C. Cyr
Computational Mathematics Department
Sandia National Laboratotories
eccyr@sandia.gov

Paul Lin
Lawrence Berkeley National Laboratory
paullin@lbl.gov

**MS40**

**Efficiently Implementing FE Boundary Conditions Using Stream-Orchestrated Execution on GPU**

The Finite Element method is commonly used to solve PDEs. The assembly procedure can easily be made massively parallel because of the natural parallelism it exposes. Frameworks like *Kokkos* can help write a single performance portable assembly code. The assembly can be expected to maximise GPU usage because the number of work units involved is typically much larger than the available concurrency. However, after the assembly, boundary conditions are oftentimes applied on small work sets, thereby potentially under-utilising the GPU processors. Moreover, boundary conditions can showcase dependencies that must be observed. In this talk, we focus on how we used GPU streams for concurrent execution of boundary condition functors in order to maximise GPU occupancy while observing their dependencies. More specifically, we will show how we used *Kokkos::Graph* to achieve such a goal in a performance portable way. We will also discuss how we designed a polymorphic hierarchy of functors for applying boundary conditions (Dirichlet, Neumann, periodic, ...) and how we map the device polymorphic functors to nodes in the graph. The design and its performance (occupancy, impact of polymorphism, ...) will be analysed in the context of electromagnetism applications.

Romin Tomasetti, Maarten Arnst
Universite de Liege

romin.tomasetti@uliege.be, maarten.arnst@uliege.be

## MS40

### A Retrospective on Performance Portable Finite Element Assembly in Albany

Albany is an open-source, C++, finite element code base written for the solution and analysis of multiphysics problems. The finite element approach implemented in Albany is designed to easily incorporate multiple physics models with graph-based evaluation and embedded analysis using template-based generic programming. Finite element assembly is decomposed into a set of nodes called "evaluators" where the performance portable framework, Kokkos, is used to support on-node, parallel execution on shared memory architectures (e.g., GPUs). Ten years have passed since Kokkos was first introduced into Albany. This talk looks back on successes and challenges of developing and maintaining performance portable finite element assembly in Albany.

Jerry Watkins, Max Carlson, Irina K. Tezaur
Sandia National Laboratories
jwatkin@sandia.gov, maxcarl@sandia.gov, ikalash@sandia.gov

## MS41

### Good Practices for Scientific Machine Learning

Scientific machine learning is a rapidly growing field, combining the predictive power of machine learning algorithms with the rigor of the physical sciences. A particular challenge of using machine learning in scientific settings is providing sufficient interpretability and transparency of the models. Another challenge in the field is reproducibility of results, which suffers at the demands to publish quickly in a rapidly evolving field, together with a lack of community standards. As with all new modeling approaches, the community benefits from increasing awareness of good practices for the development, reporting and critical assessment of the models. This talk will discuss such practices, including: identifying and quantifying all sources of uncertainty, full reporting of computational costs, model selection and hyper parameter tuning, verification and validation analysis, describing the domain of applicability and limitations, and comparing outcomes with existing technologies. The aim is to kick-start a longer conversation in the field, leading to consensus-based good practices for scientific machine learning.

Lorena A. Barba
Department of Mechanical and Aerospace Engineering
George Washington University
labarba@gwu.edu

## MS41

### Using Biased Gradients to Achieve Parallelism in Neural Network Training

Training of deep neural networks is computationally costly due to slow convergence of optimizers like stochastic gradient descent, coupled with limited parallelism in the forward and backward propagation dimension in certain dimensions. The later point implies that there is an upper bound on the strong scaling of these methods using readily available data and model parallelism. Recent work has suggested that neural networks can be trained using approximate gradient computations, thus opening a potential to increased parallelism. One technique, layer-parallel training, uses a multigrid-in-time approach to accelerate training when the depth of the neural network is very large. This talk examine's a theoretical justification for the training success of these methods. Mainly, under appropriate learning rate schedules, gradients that are statistically biased can be used in stochastic gradient descent algorithms to train deep learning models. This talk will describe the theory, while also presenting supporting numerical results.

Eric C. Cyr
Computational Mathematics Department
Sandia National Laboratotories
eccyr@sandia.gov

Matthias Heinkenschloss
Computational Applied Mathematics and Operations Research
Rice University
heinken@rice.edu

Shengchao Lin
Rice University
shengchao.lin@rice.edu

## MS41

### The Effect of Uncertainty on Learning Dynamics of Neural Networks

A typical learning process of a neural network can be described by the interaction between the data, the weights and the model architecture. In a standard regime, this interaction involves updating the weights of the neural network model with respect to a series of data batches for a fixed architecture. Multitude of decisions such as the number of layers, the activation function, learning rate, etc; are made during this process of learning. As these decisions are made apriori in the standard setup, there is an inherent uncertainty in the learning process due to the discrepancy between these decisions and the nature of the data. In this talk, we will demonstrate that this uncertainty could be quantified and further leveraged to improve performance in a neural network training process. Towards this end, we will enunciate a novel mathematical formulation where the learning problem is formulated as a dynamical system and the dynamics are represented by a differential equation. Subsequently, we will elucidate, "how to quantify this uncertainty during learning?" and provide insights into the effect of the this uncertainty. We will end this talk with insights on "how to use uncertainty to correct the behavior of the neural network model while learning?"

Krishnan Raghavan
Mathematics and Computer Science
Argonne National Laboratory
kraghavan@anl.gov

Vishwas Rao

Argonne National Laboratory
vhebbur@anl.gov

## MS42
### The Exascale Computing Project

The duration and scale of the of the Exascale Computing Project (ECP) provided a unique opportunity to advance computational science in a wide variety of application areas. The project lasted seven years and funded over one thousand researchers to develop 24 applications, over 70 software products, and deploy them on the exascale computers at DOE leadership computing facilities. This work has resulted in many notable outcomes including: the first integrated HPC software stack comprising over 100 libraries, performance portable applications refactored for accelerator-based computing architectures, and a new generation of computational scientists exposed to state-of-the-art techniques. In this talk we will give an overview of the Exascale Computing Project and highlight key legacy outcomes that will impact computational science community for years to come.

Lori A. Diachin
Lawrence Livermore National Laboratory
diachin2@llnl.gov

## MS42
### Collaborations Between ECP and the DOE HPC Facilities

The Exascale Computing Project (ECP) built an ecosystem that nourished the process of using HPC for scientific discovery from inception to execution. A key link in this chain was enabled by collaborations with the DOE HPC facilities to deploy exascale software (SD) and applications (AppInt) on the machines as early as possible, which helped debug both the machines and the software. Another key development was a Continuous Integration capability that could function in this environment. AppInt and SD teams were critical to the successful execution of exascale applications on first-of-their kind hardware in a software environment that was rapidly developing and stabilizing. We will discuss this process and give some lessons learned that could be applicable beyond ECP.

Richard Gerber
Lawrence Berkeley National Laboratory
ragerber@lbl.gov

Susan Coghlan
Argonne National Laboratory
smc@alcf.anl.gov

## MS42
### Delivering Reusable RD Software as a Construction Project: Strategies and Lessons Learned

The Exascale Computing Project (ECP) was funded and executed as a tailored construction project using earned-value management concepts, techniques, and workflows. In this presentation, we discuss the approach ECP used for library and tool software to plan, execute, track, and assess progress while still managing emerging requirements, uncertain timelines, and the co-delivery of the computing platforms for which we were targeting our work. We discuss the outcomes from this effort: a multi-layered, multi-component software stack that works portably across diverse high concurrency computing systems, containing libraries and tools that were co-designed and developed with application users and delivered as part of a robust, curated software portfolio. We discuss lesson learned from these efforts and how we intend to carry these lessons forward in future scientific software efforts.

Lois C. Mcinnes
Mathematics and Computer Science Division
Argonne National Laboratory
curfman@anl.gov

Michael A. Heroux
Sandia National Laboratories
St. John's University
maherou@sandia.gov

## MS42
### The Impact of ECP Applications

The Exascale Compute Projects Application Development portfolio, a diverse set of 24 application projects and 6 co-design centers, provides a robust demonstration of the challenges and opportunities afforded by exascale computing. ECP facilitated an unprecedented level of coordinated development across discipline boundaries, a practice which was ultimately critical to the success of the project. In this talk, we present observations on how cross-project collaboration and engagement was used within the ECP, focusing on both the benefits and challenges of greater interdependence among loosely connected communities. We share common themes that emerged, and offer advice to application teams looking to build on ECPs initial success.

Erik W. Draeger
Lawrence Livermore Nat. lab.
draeger1@llnl.gov

Andrew Siegel
Argonne National Laboratory
siegela@uchicago.edu

## MS43
### Accelerating Sparse Solvers with Cache-Optimized Matrix Power Kernels

Sparse linear iterative solvers are indispensable for large-scale simulations. In this talk, we present methods to accelerate some of the existing solvers and preconditioners by using the concept of levels as developed in the context of our RACE library framework. Levels are constructed using breadth-first search on the graph related to the underlying sparse matrix. These levels are then used to implement cache blocking of the matrix elements for high spatial and temporal reuse. The approach finds its use in kernels like sparse-matrix-power vector multiplication, which perform repetitive back-to-back sparse-matrix-vector multiplication (SpMV) operations. The method is highly effective

and achieves performance levels of 50-100 GF/s on a single modern Intel or AMD multicore chip, providing speedups of typically 2x-4x compared to a highly optimized classical SpMV implementation. After briefly introducing the optimization strategy, we shed light on the application of these optimized kernels in iterative solvers. To this end, we discuss the coupling of the RACE library with the Trilinos framework and address the application to communication-avoiding s-step Krylov solvers, polynomial preconditioners, and algebraic multigrid preconditioners. We then dive into the performance benefits and challenges of the RACE integration and show that our optimization produces numerically identical results and improves the total solver time by 1.3x-2x.

Christie Louis Alappat
Friedrich-Alexander-Universität Erlangen-Nürnberg
Erlangen, Germany
christie.alappat@fau.de

## MS43

### Sparse Algorithms for Large-Scale Bayesian Inference Problems

Bayesian inference provides a mean to predict events or uncover dependency structures within complex systems while simultaneously quantifying the uncertainty of its own estimates. It is widely used with applications across many disciplines including life sciences, medicine, and geostatistics. The underlying statistical models can often be formulated using sparse Gaussian Markov random fields. They give rise to matrices with recurring sparsity patterns allowing for the usage of sparse solution methods. Foremost ones requires efficient algorithms for Cholesky decomposition and selected matrix inversion. We give an overview of different approaches, ranging from entirely sparse methods to block solvers. In particular we compare double, single and mixed precision algorithms in this context, present their performance and the effects on the quality of the results.

Lisa Gaedke-Merzhäuser, Olaf Schenk
Università della Svizzera italiana
lisa.gaedke.merzhaeuser@usi.ch, olaf.schenk@usi.ch

## MS43

### Efficient Schwarz Preconditioning Techniques on Current Hardware Using FROSch

FROSch (Fast and Robust Overlapping Schwarz) is a framework for parallel Schwarz domain decomposition preconditioners in Trilinos. It is applicable to a wide range of problems due to its algebraic approach, which allows the preconditioners to be constructed from a fully assembled, parallel distributed matrix. This is enabled by the use of extension-based coarse spaces instead of classical coarse spaces. FROSch also features a variety of algorithmic variants that extend its applicability and scalability. These include monolithic preconditioning for block systems and multi-level extensions. This talk will focus on recent developments in FROSch that improve the efficiency of Schwarz preconditioners for current hardware architectures. These include techniques for reducing communication and global work, lower precision preconditioning, as well as techniques

for facilitating GPUs. The performance of the different techniques will be demonstrated for different application problems and using different state-of-the-art supercomputers.

Alexander Heinlein
Delft University of Technology
Delft Institute of Applied Mathematics
a.heinlein@tudelft.nl

Sivasankaran Rajamanickam, Ichitaro Yamazaki
Sandia National Laboratories
srajama@sandia.gov, iyamaza@sandia.gov

## MS43

### Communication-Reduced Sparse-Dense Matrix Multiplication with Adaptive Parallelization

This work presents the Communication-Reduced Parallel SpMM (CRP-SpMM) algorithm, a simple and novel algorithm that adaptively parallelizes SpMM to reduce communication costs. Unlike existing 1.5D and 2D distributed-memory parallel SpMM algorithms that use the same matrix and workload partitioning for the same number of processes, CRP-SpMM selects the matrix partitionings and parallelization schemes based on the sparsity pattern as well as the dimensions of the input matrices. This allows CRP-SpMM to generalize 1D and 2D algorithms and choose a better parallelization approach for the problem than existing approaches. CRP-SpMM uses a two-phase communication scheme to improve performance on the required communication operations. Numerical experiments show that CRP-SpMM significantly outperforms existing distributed-memory parallel SpMM algorithms.

Hua Huang
Georgia Institute of Technology
huangh223@gatech.edu

Edmond Chow
School of Computational Science and Engineering
Georgia Institute of Technology
echow@cc.gatech.edu

## MS44

### Unified Data Abstractions for Scientific Workflow Composition in the Computing Continuum

Nowadays, complex scientific workflows must be able to extract knowledge and produce insight at every step from the instruments to the scientist. This requires the integration of heterogeneous and geographically distributed computing environments into a computing continuum to seamlessly bridge simulations, machine learning and data-driven analytics. Current approaches to large-scale computing are accommodating the need to support hyper-heterogeneous environments in which different ecosystems and devices work together. Existing works on workflow composition and deployment in the computing continuum focus on task-flow control and are disconnected from data patterns and structures beyond domain-specific applications. Moreover, general approaches for representing knowledge and provenance in the form of metadata are also lacking for these converged

workflows, and common interfaces for data management in the continuum are necessary. Unified data abstractions can enable the interoperability of data storage and processing across the continuum and facilitate data analytics at all levels, alleviating the disconnect between application- and storage-oriented approaches to interoperability. This talk presents use cases and new research directions towards unified data modeling approaches to structure and represent data on a logical level across the computing continuum.

Silvina Caino-Lores
Inria Rennes Bretagne Atlantique
scaino-lores@acm.org

## MS44

### Data Movement and the Data Bridge in the Destination Earth Climate Adaptation Digital Twin

The rise of Digital Twins to help design and test disruptive technology spawns virtual laboratories where data streams between real life measurements and computational models. Destination Earth, a European Flagship Initiative to develop a highly accurate digital model of the Earth on a global scale, is one such example. It will be utilizing the EuroHPC Pre-Exascale and Exascale systems, federating workload and potentially workflows. Such considerable influx of data exacerbates the well-known High Performance Computing (HPC) problem of dealing with data-intensive workloads. Deficiencies of systems software appear at multiple levels of the stack, be it as high as the workflow management level in orchestrating and coupling applications, or as low as the memory hierarchy level where memory and storage technologies converge carrying incompatible software in managing data. Furthermore, data governance and protection topics need to be addressed inside each workflow, not just in the data acquisition and product dissemination phase. We will discuss the architecture and challenges, and project status at the end of the first project phase, from the perspective of the Climate Adaptation twin.

Utz-Uwe Haus
EMEA Research Lab, Hewlett Packard Enterprise
utz-uwe.haus@hpe.com

Ali Mohamed
, EMEA Research Lab, Hewlett Packard Enterprise
mohammed, ali ¡ali-omar-abdelazim.mohammed@hpe.com

Christopher Haine
EMEA Research Lab, Hewlett Packard Enterprise
haine, christopher ¡christopher.haine@hpe.com¿

## MS44

### Challenges in Privacy-preserving Federated Learning Workflows on DOE HPC Resources

Federated learning (FL) is a collaborative learning approach where multiple data owners, referred to as clients, train a model together under the orchestration of a central server by sharing the model trained on their local datasets instead of sharing the data directly. FL enables creation of more robust models without the exposure of local datasets.

However, FL by itself, does not guarantee the privacy of data, because the information extracted from the communication of FL algorithms can be accumulated and utilized to infer the private local data used for training. We developed the Argonne Privacy Preserving Federated Learning framework to enable Privacy-Preserving Federated Learning. We enabled training of AI models in a distributed setting across multiple institutions, where sensitive data are located, with the ability to scale on supercomputing resources to help create robust, trust-worthy AI models in biomedicine and smart grid applications where data privacy is essential. Setting up a secure federated learning experiment that needs HPC resources across distributed sites requires technical capabilities that may not be available for all. To lower the barrier to entry for leveraging PPFL and to enable domain experts in large institutions to utilize FL, we created the Argonne Privacy-Preserving Federated Learning as a service (APPFLx), which enables cross-silo PPFL using easy to use web interface for managing, deploying, analyzing, and visualizing PPFL experiments.

Ravi Madduri
Argonne National Laboratory
madduri@anl.gov

## MS44

### Challenges and Requirements to Drive Workflows from the Data Plane

We observe the emergence of a new generation of scientific workflows in which the efficient management of the data produced and consumed by the different workflow components become more important than the efficient execution of the computational part. In this talk, we identify several challenges that these workflows cause to traditional workflow systems based on the analysis of five generic workflow motifs originating from diverse scientific domains. We also express some requirements that a modern workflow ecosystem should meet to efficiently drive workflows for their data plane.

Frederic Suter
Oak Ridge National Laboratory
suterf@ornl.gov

## MS45

### Coupling Earthquake Cycles with Plate Tectonics: Toward Scalable Numerical Methods for Multiscale Visco-Elastic-Plastic Pdes

The Rhea mantle convection code, previously developed for solving the Stokes equations with a visco-plastic (VP) constitutive relation, has demonstrated scalability on hundreds of thousands of cores by using aggressive adaptive mesh refinement on hexahedral meshes, accurate and efficient high-order tensor-product finite elements, and a highly advanced Newton-Krylov solver for the resulting nonlinear system of equations. We now propose to apply the Rhea code to a visco-elastic-plastic (VEP) Stokes system, where the elasticity makes the constitutive relation time dependent and the plasticity is interpreted as a complementarity condition involving the limiting time-

dependent yield stress. As we will demonstrate, implicit temporal discretization of the constitutive relation leads to a VP-type Stokes system which must be solved at each timestep, allowing direct use of the solvers from Rhea. We provide a roadmap toward achieving our major application goal of attaining realistic simulations coupling plate tectonics with earthquake cycles. We also discuss our progress in overcoming new computational challenges, such as highly localized features in time and space that appear during rupture events, and present preliminary results from model problems implemented in a toy Firedrake code and Rhea.

Max P. Heldman
Boston University
Argonne National Laboratory
maxh@vt.edu

Georg Stadler
Courant Institute for Mathematical Sciences
New York University
stadler@cims.nyu.edu

Michael Gurnis
Seismological Laboratory
California Institute of Technology
gurnis@gps.caltech.edu

Johann Rudi
Virginia Tech
jrudi@vt.edu

Jiaqi Fang
Seismological Laboratory
California Institute of Technology
jfang@caltech.edu

## MS46

### FEAST Parallel Eigenvalue Solver with Applications to First-Principle Calculations

Realistic first-principle quantum simulations applied to large-scale atomistic systems pose unique challenges in the design of eigenvalue algorithms that are both capable of processing a considerable amount of generated data, and achieving significant parallel scalability on modern high-end computing architectures. The FEAST eigensolver [www.feast-solver.org] represents a unified framework for solving linear and non-linear eigenvalue problems and achieving accuracy, robustness, and high-performance scalability on parallel architectures. FEAST offers a large number of features including three-levels of MPI parallelisms, mixed-precision arithmetic, and iterative or direct factorization approach. From atoms and molecules to nanostructures, we discuss how the FEAST eigensolver can optimally be used for enabling reliable and high-performance large-scale first-principle calculations. We will present various simulation results obtained using our all-electron real-space DFT and real-time TDDFT NESSIE software [www.nessie-code.org].

Dongming Li
University of Massachusetts, Amherst, U.S.
dongmingli@umass.edu

Eric Polizzi
University of Massachusetts, Amherst, USA
epolizzi@engin.umass.edu

## MS46

### Large Scale Low Rank Approximations: Streaming vs Randomized Sketching

With the unprecedented growth in size, data often cannot be accessed more than once and sometimes not even stored. In such cases, dimensionality reduction or low rank approximation is often performed through randomized sketching methods. Streaming methods provide an alternative where one or a small number of matrix rows are read at a time and their contribution is accounted for before they are discarded. We consider a generalization of this idea, where the block of rows stored at a time is as large as can be stored, and the update occurs through iterative methods. We study the trade-offs of streaming over randomized sketching methods in terms of accuracy, complexity, and parallelization, especially from the viewpoint of large scale rather than online applications.

Andreas Stathopoulos
College of William & Mary
Department of Computer Science
andreas@cs.wm.edu

Jeremy Myers
The College of WIlliam and Mary
jmmyers01@email.wm

## MS46

### Parallel Orthogonality in Large-Scale Tensor Network Eigenvalue Problems

Orthogonality plays a key role in eigenvalue computations. In 1D tensor networks such as tensor trains, the orthogonality is maintained by using QR or truncated SVD factorizations. However, this technique does not extend to 2D tensor networks such as projected entangled pair states (PEPS). Moreover, orthogonality inside a PEPS keeps the computational complexity of eigenvalue evaluations bounded. We will discuss and compare several approximate orthogonalization techniques and strategies for orthogonalizing PEPS columns and rows.

Roel Van Beeumen
Lawrence Berkeley National Laboratory
rvanbeeumen@lbl.gov

## MS46

### Advancing Chase Library Towards Exascale Applications on Distributed Multi-GPUs and ARM-based Systems

With the continuous evolution of supercomputers, characterized by their expanding size and the integration of powerful Graphics Processing Units (GPUs), traditional direct eigensolvers are confronted with significant challenges in keeping pace with this hardware transformation. The primary hurdles they face stem from the increasing demands for communication and synchronization, hindering their

scalability. In contrast, subspace eigensolvers offer a more streamlined structure that enables them to surmount these bottlenecks effectively. ChASE represents a contemporary subspace eigensolver renowned for its utilization of Chebyshev polynomials, enhancing the computation of extremal eigenpairs of dense Hermitian eigenproblems. This study showcases a series of innovative modifications applied to ChASE, including a thorough reevaluation of its memory layout, the introduction of a mixed 1D/2D parallelization approach, the adoption of a more efficient communication-avoiding algorithm for one of its core modules, and the substitution of the MPI library with the vendor-optimized NCCL library when used on NVIDIA GPUs. The outcome of these enhancements is a library capable of seamlessly handling dense problems of considerable scale, reaching dimensions up to $N = \mathcal{O}(10^6)$ on accelerated clusters and ARM-based platforms.

Edoardo A. Di Napoli
Juelich Supercomputing Centre
e.di.napoli@fz-juelich.de

Xinzhe Wu
Jülich Supercomputing Centre
xin.wu@fz-juelich.de

**MS47**

**Iterative Linear Solvers in Interior Point Methods for Large-Scale Optimization**

The solution of large, sparse and structured matrices lies at the heart of many state-of-the-art algorithms for solving optimization problems. As an example, interior point methods, which are widely used for nonlinear and linear programming alike, rely on (a variant of) Newton's method to find search directions at each iteration. The matrices that arise in interior point methods become extremely ill-conditioned by construction as an optimal solution is approached. This ill-conditioning occurs in a very specific manner, and has been shown to be benign, both from a more theoretical standpoint and in practice. Many implementations of interior point methods available use direct methods for solving linear systems, which has proven to work well in spite of the aforementioned ill-conditioning. However, the need to solve ever larger optimization problems and the rise of hardware accelerators in computing presents a challenge, since The topic of this talk will be to discuss how iterative linear solvers may be used for interior point methods instead, with key challenges lying in the design of suitable preconditioners and handling of the poor conditioning of the systems. We will attempt to both give a background on the challenges in the area and some of our recent work on the topic.

Felix Liu
Raysearch Laboratories
KTH Royal Institute of Technology
felixliu@kth.se

**MS47**

**Give Us Cache, We Give You Bandwidth!**

This talk presents numerical algorithms recent impacts on solving scientific problems using massively parallel systems. Based on algebraic compression, we illustrate the algorithmic approach on seismic imaging, computational astronomy, and climate modeling using x86 and accelerator-based systems. We also assess the impact on energy efficiency and identify the need for cross-disciplinary expertise to further address one of the most urgent challenges faced by the scientific community.

Hatem Ltaief
KAUST, Saudia Arabia
hatem.ltaief@kaust.edu.sa

**MS47**

**Batched Sparse Solvers, Libraries, Preconditioners, and BLAS-like Functionality for Modern Platforms**

Many applications utilizing Exascale hardware with GPU-based nodes rely on computations on many small matrices, that are often sparse. The common operation is to obtain simultaneous solutions of multiple linear systems of equations that are structurally sparse. Consequently, batched sparse linear solvers form the new frontier for algorithmic development and require a new toolset for their performance engineering. With high efficiency in mind, we strive to design, implement, and integrate into these new breed of applications, that require the performance levels available from modern systems, which tightly combine accelerators, high-bandwidth memory, and network cards. Providing appropriate interface will be described, which we believe enables, rather than inhibits, efficient implementation and the use of batched sparse functionality from solvers and preconditioners to other BLAS-like routines. In addition to the interface design considerations, we will present existing implementations under active development and some of the results for the relevant set of application domains.

Piotr Luszczek
Massachusetts Institute of Technology
luszczek@ll.mit.edu

Stanimire Tomov
University of Tennessee
tomov@icl.utk.edu

**MS47**

**Residual Simplex Method**

We propose a novel method for solving large-scale linear programming problems based on ideas from Krylov subspace methods. Each iteration we solve a projected problem where the linear constraints are projected on the basis of the residuals from the previous iterations. This results in a series for small projected LP problems that increase in size. But each problem can be warm-started with the basic set and of solution of the previous problem. We discuss the convergence, the parallel implementation and the performance on some large-scale network problems.

Wim I. Vanroose, Bas Symoens
University of Antwerpen
wim.vanroose@uantwerpen.be,

bas.symoens@uantwerpen.be

## MS48

**Tools to Capture Message Passing Characteristics of HPC Applications and Support Design of Realistic HPC Benchmarks for Heterogeneous Systems**

The current set of proxy and mini applications used to benchmark heterogeneous systems lack realistic datasets that can be used to evaluate the performance of the actual applications they intend to represent. To address this issue of unavailability of necessary realistic datasets, in this talk we explore tools that could be used to capture the message passing characteristics of the actual application executing with real datasets on current HPC systems and then use this information to create realistic HPC benchmarks. This approach will ensure that the characteristics of the actual application are captured without requiring the application developers to modify the existing application as well as not share the datasets used to execute the application. Only details required to understand the message passing characteristics of the application are captured and then used to create benchmarks that are representative of the actual applications.

Purushotham Bangalore
University of Alabama
pvbangalore@ua.edu

## MS48

**OpenHPCA: High Performance Compute Availability Benchmark for Smart Networks**

The High Performance Compute Availability (HPCA) Benchmark project is an effort to create a new metric for ranking HPC and AI system performance and capabilities. HPCA is intended as a complement to existing benchmarks, such as HPCC and HPCG. NVIDIA, as co-chair of the group, led the design, integration, and implementation of a comprehensive set of benchmarks to evaluate the overall compute resource performance in the presence of in-network computing technologies, e.g., NVIDIA BlueField hardware. These benchmarks are a mix of existing benchmarks as well as new benchmarks. In this context, we will present the latest version of our benchmark suite, which is a combination of existing and new micro-benchmarks. We will show how it can be used to quantify the benefit of using smart network technologies, e.g., by offloading MPI operations to BlueFields cards. Finally, we will present an overview of the current directions for new benchmarks discussed by the group.

Richard Graham, Geoffroy Valle
NVIDIA
richgraham@nvidia.com, geoffroyvalle@nvidia.com

## MS48

**A Strategy and Approach for Generating Synthetic Message-Passing Benchmark Data from Real Applications**

Many modern scientific applications have communication patterns in which both the number of communication part-

ners and amount of data transmitted between process pairs vary significantly plus change over time. This talk describes an approach that provides a lightweight method to reproduce communication patterns of a variety of applications with minimal overhead while gaining additional insights into the performance and characteristics of various irregular communication patterns. First, an approach to measure and model the behavior of these irregular, dynamic MPI communication patterns on modern high-performance computing systems is described. This approach quantifies communication behavior using a small number of stochastic random variables that capture key features of irregular communication patterns, and estimates the distributions of these variables either parametrically or empirically. Overall, this work demonstrates that the collected parameters and their distributions can be used to measure and model the communication performance of several MPI applications. We present a synthetic benchmark that uses these distributions to recreate statistically similar communication patterns.

Patrick G. Bridges
University of New Mexico
bridges@cs.unm.edu

Carson Woods
Emory University
carson.woods@emory.edu

Derek Schafer
University of New Mexico
dschafer1@umn.edu

Anthony Skjellum
University of Tennessee at Chattanooga
askjellum@tntech.edu

## MS48

**The Impact of Space-Filling Curves on Data Movement in Parallel Systems**

Modern computer systems are characterised by deep memory hierarchies, composed of main memory, multiple layers of cache, and other specialised types of memory. In parallel and distributed systems, additional memory layers are added to this hierarchy. Achieving good performance for computational science applications, in terms of execution time, depends on the efficient use of this diverse and hierarchical memory. Furthermore, the pertinent trade-offs change as application, memory, and processor characteristics evolve. With these changes in mind, this paper revisits the use of space-filling curves to specify the ordering in memory of data structures used in representative scientific applications executing on parallel machines containing clusters of multicore CPUs with attached GPUs. This work examines the hypothesis that space-filling curves, such as Hilbert and Morton ordering, can improve data locality and hence result in more efficient data movement than row or column-based orderings. First, performance results are presented that show for what application parameterisations and machine characteristics this is the case, and are interpreted in terms of how an application interacts with the computer hardware and low-level software. This

research particularly focuses on the use of stencil-based applications that form the basis of many scientific computations. Second, how space-filling curves impact data sharing in nearest-neighbor and stencil-based codes is considered.

David Walker
Tennessee Technical University
david.walker611@gmail.com

Anthony Skjellum
University of Tennessee at Chattanooga
askjellum@tntech.edu

## MS49

### Runtime Systems for Scalable Hierarchical Matrix Factorization Algorithms

Low rank approximation of structured dense matrices can reduce the computation and storage requirements of dense matrix factorization. Low rank approximation can be done with the help of shared and nested basis using the $\mathcal{H}^2$-matrix format. The $\mathcal{H}^2$-ULV factorization can be used for $O(N)$ factorization of the $\mathcal{H}^2$-matrix. Similar to the dense LU factorization, the $\mathcal{H}^2$-ULV factorization has off-diagonal dense dependencies. However, the number of dense blocks is much lesser than that of the dense factorization. Although variants of the $\mathcal{H}^2$-ULV without off-diagonal dependencies have been proposed, they have been shown to be inaccurate for ill-conditioned matrices. In this talk, we will see how the use of runtime systems can improve the efficiency of $\mathcal{H}^2$-ULV factorization with dependencies on distributed supercomputers such as Fugaku. We demonstrate the use of runtime systems for the $\mathcal{H}^2$-ULV factorization for algorithms such as direct solvers and for the slicing the spectrum algorithm for obtaining eigen values.

Sameer Deshmukh
Tokyo Institute of Technology
sameer.deshmukh@rio.gsic.titech.ac.jp

Qianxiang Ma
Dept. of Computer Science
Tokyo Institute of Technology
ma@rio.gsic.titech.ac.jp

Rio Yokota
Tokyo Institute of Technology
rioyokota@gsic.titech.ac.jp

George Bosilca
University of Tennessee Knoxville
bosilca@icl.utk.edu

Ridwan Muhammed
Dept. of Computer Science
Tokyo Institute of Technology
ridwan@rio.gsic.titech.ac.jp

## MS49

### Scaling the Memory Wall with Tile Low Rank

We exploit the high memory bandwidth of AI-customized Cerebras CS-2 systems for seismic processing. Through low-rank matrix approximation, memory hungry seismic applications fit onto memory-austere SRAM waferscale hardware, addressing a challenge arising in many wave-equation-based algorithms that rely on multi-dimensional convolution operators. Exploiting sparsity inherent in seismic data in the frequency domain, we implement embarrassingly parallel tile low-rank matrix-vector multiplications (TLR-MVM), which account for most of the elapsed time in MDC operations, to solve the Multi-Dimensional Deconvolution (MDD) inverse problem. By reducing memory footprint along with arithmetic complexity, we fit a standard seismic benchmark dataset into the local memories of Cerebras processing elements. TLR-MVM on 48 CS-2 systems in support of MDD gives a sustained memory bandwidth of 92.58PB/s on 35,784,000 processing elements, a significant milestone that highlights the capabilities of AI-customized architectures to enable a new generation of seismic algorithms that will empower multiple technologies of our low-carbon future. We compare this tile low-rank computation, which was a Gordon Bell Prize finalist for 2023, with another tile low-rank computation based on a Cholesky factorization of a covariance matrix from environmental statistics run on Fugaku, which was a Gordon Bell Prize finalist for 2022.

David E. Keyes
King Abdullah University of Science and Technology (KAUST)
david.keyes@kaust.edu.sa

Hatem Ltaief
KAUST, Saudia Arabia
hatem.ltaief@kaust.edu.sa

Yuxi Hong, Matteo Ravasi
King Abdullah University of Science and Technology
yuxi.hong@kaust.edu.sa, matteo.ravasi@kaust.edu.sa

Mathias Jacquelin, Leighton Wilson
Cerebras
mathias.jacquelin@cerebras.net,
leighton.wilson@cerebras.net

## MS49

### Construction of Hierarchically Semi-Separable Matrix Using Faster Randomized Sketching

We extend our early work on adaptive partially matrix-free Hierarchically Semi-Separable (HSS) matrix construction algorithm using Gaussian sketching operators to a broader class of Johnson-Lindenstrauss (JL) sketching operators. We present theoretical work which justifies this extension. In particular, we extend the earlier concentration bounds to all JL sketching operators and examine this bound for specific classes of such operators including the original Gaussian sketching operators, subsampled randomized Hadamard transform (SRHT) and the sparse Johnson-Lindenstrauss transform (SJLT). We demonstrate experimentally that using SJLT instead of Gaussian sketching operators leads to 1.5–2.5× speedups of the HSS construction implementation in the STRUMPACK C++ library. The generalized algorithm allows users to select

their own JL sketching operators with theoretical lower bounds on the size of the operators which may lead to faster runtime with similar HSS construction accuracy.

Xiaoye S. Li
Computational Research Division
Lawrence Berkeley National Laboratory
xsli@lbl.gov

Yotam Yaniv
Univ. of California Los Angeles
yotamya@math.ucla.edu

Osman Malik
Lawrence Berkeley National Laboratory
oamalik@lbl.gov

Pieter Ghysels
Lawrence Berkeley National Laboratory
Computational Research Division
pghysels@lbl.gov

## MS49

### Butterfly Algorithms: from Matrix to Tensor

Butterfly algorithms represent a rapidly growing class of multi-level low-rank decomposition algorithms, which are not only well-suited for compressing highly oscillatory operators including special integral transforms and Green's functions for wave equations, but also lately applied to high-dimensional PDEs and machine learning. This talk summarizes some of the recent development in butterfly algorithms including: 1. New development in (a) sketching, (b) arbitrary entry-evaluation and (c) completion-based fast construction algorithms. 2. New strongly-admissible hierarchical matrix-based algorithm that achieves same complexity as the fast multipole methods. 3. Tensor extension of matrix butterfly algorithm which outperforms existing tensor algorithms such as tensor-train or Tucker for high-dimensional oscillatory operators.

Yang Liu
Lawrence Berkeley National Laboratory
liuyangzhuan@lbl.gov

## MS50

### Establishing Fundamental Lower Bounds for Algorithmic Strategies Using Time-dependent Hamiltonians

Many hybrid and near-term quantum algorithms rely on time-dependent Hamiltonians. In particular, at each iteration i of an algorithm, one may prepare a quantum state | ??? through the application of a time-varying Hamiltonian H(t) to initial state | ???1?. In this talk, I will introduce and prove a computationally-localized adiabatic theorem, which allows one to analyze such strategies as inherently adiabatic processes within the localized framework. This facilitates the establishment of concrete lower bounds on the efficiency and performance of these approaches. This talk is based upon joint work with Jacob Bringewatt and

Connor Mooney.

Michael Jarrett Baume
George Mason University
mjarretb@gmu.edu

Jacob Bringewatt, Connor Mooney
University of Maryland
jbringew@umd.edu, tmooney@umd.edu

## MS50

### Fast Learning To Optimize Quantum Neural Network Without Gradients

Quantum Machine Learning is an emerging sub-field in machine learning where one of the goals is to perform pattern recognition tasks by encoding data into quantum states. This extension from classical to quantum domain has been made possible due to the development of hybrid quantum-classical algorithms that allow a parameterized quantum circuit to be optimized using gradient based algorithms that run on a classical computer. The similarities in training of these hybrid algorithms and classical neural networks has further led to the development of Quantum Neural Networks (QNNs). However, in the current training regime for QNNs, the gradients w.r.t objective function have to be computed on the quantum device. This computation is highly non-scalable and is affected by hardware and sampling noise present in the current generation of quantum hardware. In this paper, we propose a training algorithm that does not rely on gradient information. Specifically, we introduce a novel meta-optimization algorithm that trains a meta-optimizer network to output parameters for the quantum circuit such that the objective function is minimized. We empirically and theoretically show that our method is competitive with existing gradient based algorithms and outperforms them in terms of computation time on different datasets

Ankit Kulshrestha
University of Delaware
akulshr@udel.edu

Xiaoyuan Liu
Fujitsu Research USA
xliu@fujitsu.com

Hayato Ushijima-Mwesigwa
Fujitsu Corp, Japan
hayato@us.fujitsu.com

Ilya Safro
University of Delaware
isafro@udel.edu

## MS50

### Of Representation Theory and Quantum Approximate Optimization Algorithm

In this paper, the Quantum Approximate Optimization Algorithm (QAOA) is analyzed by leveraging symmetries inherent in problem Hamiltonians. We focus on the generalized formulation of optimization problems defined on

the sets of n-element d-ary strings. Our main contribution encompasses dimension reductions for the originally proposed QAOA. These reductions retain the same problem Hamiltonian as the original QAOA but differ in terms of their mixer Hamiltonian, and initial state. The vast QAOA space has a daunting dimension of exponential scaling in $n$, where certain reduced QAOA spaces exhibit dimensions governed by polynomial functions. This phenomenon is illustrated in this paper, by providing partitions corresponding to polynomial dimensions in the corresponding subspaces. As a result, each reduced QAOA partition encapsulates unique classical solutions absent in others, allowing us to establish a lower bound on the number of solutions to the initial optimization problem. Our novel approach opens promising practical advantages in accelerating the class of QAOA approaches, both quantum-based and classical simulation of circuits, as well as a potential tool to cope with barren plateaus problem.

Boris Tsvelikhovskiy
tbd
BDT18@pitt.edu

Ilya Safro
University of Delaware
isafro@udel.edu

Yuri Alexeev
Argonne National Laboratory
yuri@alcf.anl.gov

Boris Tsvelikhovskiy
tbd
BDT18@pitt.edu

## MS51

### Matrix-free Methods in deal.II using Kokkos

In recent years, matrix-free methods have proven to be favorable over matrix-based methods if the whole solver stack is compatible. The FEM framework deal.II had a CPU-only framework for at least 10 years and a Cuda-based version for about 5 years. Recently, we ported the latter to Kokkos. In this talk, we will discuss the design and implementation and discuss performance on different architectures.

Daniel Arndt, Bruno Turcksin
Oak Ridge National Laboratory
arndtd@ornl.gov, turcksinbr@ornl.gov

## MS51

### On the Performance and Portability of the FIERRO Mechanics Code

We present recent research on the performance and portability of the open-source massively parallel C++ FIERRO mechanics code. The FIERRO code offers novel capabilities to simulate multiscale mul- tiphysics engineering problems and for autonomous multiphysics design optimization. The FIERRO code contains a suite of novel Lagrangian finite element methods which have meshes with constant mass elements that move with the material to simulate quasi-static solid mechanics problems and compressible material dynamics problems. Multiscale material models are coupled to world-leading numerical methods inside the FIERRO code to simulate the performance of parts accounting for the underlying microstructure in the part. FIERRO is written with hybrid parallelism, using MPI for coarse- grained parallelism plus the Kokkos library for fine-grained parallelism on diverse GPUs and multi-core CPUs. The combined capabilities in Fierro are key to realizing the full-potential of diverse manufac- turing processes by leveraging modern high-performance computing (HPC) machines to: assess the performance of parts accounting for the underlying material microstructure, and design optimal parts autonomously that are additively manufactured. This talk will discuss the FIERRO mechanics code and present parallel scaling results for a range of HPC machines.

Nathaniel Morgan, Adrian Diaz, D Dunning, Kevin Welsh, Caleb Yensuah, Steven Walton, Sarah Hankins
Los Alamos National Laboratory
nmorgan@lanl.gov, adiaz@lanl.gov,
daniel.dunning@lanl.gov, kevin.welsh@lanl.gov,
cyenusah@lanl.gov, stevenw@lanl.gov,
sarah.hankins@lanl.gov

Svetlana Tokareva
Los Alamos National Laboratory
Applied Mathematics and Plasma Physics Group (T-5)
tokareva@lanl.gov

Evan Lieberman
LANL
elieberman3@lanl.gov

## MS51

### COO: A New Matrix Assembly Interface in PETSc to Support GPUs

PETSc is a popular math library for the scalable solution of scientific applications modeled by partial differential equations. PETSc has a rich set of application programming interfaces (APIs) for matrix assembly. PETSc users have been using them for decades in various discretization algorithms, including the finite element method. Though these APIs are powerful and efficient for assembling global matrices on CPUs at scale, they are single-thread oriented and thus not feasible for GPU computation with massive parallelism. In this talk, we will introduce our newly proposed matrix assembly interface called coordinate (COO) assembly to conquer this problem. COO allows users to provide the coordinates of the unassembled nonzeros of the matrix at once, then PETSc will figure out the sparsity pattern of the matrix and set up MPI communication plans. Later, users give PETSc the nonzero values on device in the same order as the coordinate array, and PETSc will assemble the global matrix efficiently using prebuilt plans. We implemented the COO interface with the Kokkos programming model. The interface has been used by PETSc users.

Junchao Zhang
Argonne National Laboratory
jczhang@mcs.anl.gov

Jed Brown
University of Colorado Boulder
jed@jedbrown.org

Barry Smith
Simons Foundation and Flatiron Institute, U.S.
bsmith@petsc.dev

Stefano Zampini
King Abdullah University of Science and Technology
stefano.zampini@kaust.edu.sa

## MS52

### A Comparison of Intel and OSU All-to-all Benchmarks for Next Generation FFT Algorithms

The fast Fourier transform (FFT) is an algorithm used in a wide variety of applications, including signal processing, image processing, and computational fluid dynamics. Next generation FFT algorithms are being developed to improve the performance of FFTs on parallel hardware. In this talk, we will compare the performance of the Intel and OSU all-to-all benchmarks for next generation FFT algorithms. We will focus on the startup time of the MPI_Init() function, as well as the latency and bandwidth of the all-to-all operation. We will show that the OSU benchmark consistently outperforms the Intel benchmark, especially for large numbers of processors. We will also discuss the implications of these results for optimizing the performance of next generation FFT algorithms.

Samar A. Aseeri
Computational Scientist at KAUST
samar.aseeri@kaust.edu.sa

Benson Muite
Kichakato Kizito, Kenya
benson_muite@emailplus.org

David E. Keyes
King Abdullah University of Science and Technology
(KAUST)
david.keyes@kaust.edu.sa

## MS52

### Latest Advanced on Parallel and Distributed FFT Computation on NVIDIA GPU

The present talk aim to highlight the contributions made by NVIDIA in advancing Fast Fourier Tranform (FFT) on accelerated platforms. We will cover advancements in methods, accuracy, performance across sizes, API design and multi-node scaling. Currently the NVIDIA FFT software portfolio includes cuFFT (single GPU, single node, host-side API), cufftXt (multiple GPU, single node, host-side API), cuFFTDx (cuFFT Device Extensions - single GPU device-side API), cuFFTmp (multi GPU, distributed, host-side API). All these libraries are distributes via the NVIDIA HPC SDK (https://developer.nvidia.com/hpc-sdk). NVIDIA has also developed a open-source communication library called cuDecomp. cuDecomp is a library for managing 1D (slab) and 2D (pencil) parallel decompositions of 3D Cartesian spatial domains on NVIDIA GPUs, with routines to perform global transpositions and halo communications. The library is inspired by the 2DE-COMP&FFT Fortran library, a popular decomposition library for numerical simulation codes, with a similar set of available transposition routines. The talk will cover key examples of adoption of NVIDIA libraries into other distributed libraries (e.g. 2DECOMP&FFT, P3DFFT, heFFTe) and real production codes in various domains including Quantum Chemistry (e.g. Quantum ESPRESSO, VASP) and Engineering/CFD (e.g. CaNS, Xcompact3D).

Miguel Ferrer Avila, Josh Romero, Lukasz Ligowski, Filippo Spiga
NVIDIA
mferreravila@nvidia.com, joshr@nvidia.com, lligowski@nvidia.com, fspiga@nvidia.com

## MS52

### Fftx: Release, Updates and Next Steps

We present an update on the design of the API and runtime environment of FFTX. The FFTX is developed as part of the DOE ExaScale effort by LBL, Carnegie Mellon University, and SpiralGen, Inc. We aim at translating the LA-PACK/BLAS approach from the numerical linear algebra world to the spectral algorithm domain. FFTX is extending and updating FFTW for the exascale era and beyond while providing backwards compatibility. We utilize the SPIRAL code generation system and runtime compilation as backend. A key innovation is the concept of "integrated algorithms" that allows for cross-library call optimization. This update discusses the FFTX release at the end of ECP and where the project will go next.

Franz Franchetti
Department of Electrical and Computer Engineering
Carnegie Mellon University
franzf@ece.cmu.edu

## MS52

### Implementation of Parallel Number-Theoretic Transform on GPU Clusters

In this talk, we propose an implementation of parallel number-theoretic transform (NTT) on GPU clusters. The butterfly operation of the NTT can be performed using modular addition, subtraction, and multiplication. We show that a method known as the four-step fast Fourier transform (FFT) algorithm can be applied to the NTT. We parallelized the four-step NTT using MPI and OpenACC. Performance results of parallel NTTs on GPU clusters are reported.

Daisuke Takahashi
Center for Computational Sciences
University of Tsukuba
daisuke@cs.tsukuba.ac.jp

## MS53

### Modernization Project for the Next EigenExa Library

The technology to achieve massively parallel eigenvalue

computing on supercomputers has been fixed since the success of the two-step algorithm used in the K computer or the U.S. Summit. The improvements in memory and network-bound aspects of the Householder tridiagonalization can be considered the critical factor. Furthermore, the block Householder back-transform can be accelerated by well-optimized PXGEMM routines and automatic tuning of the block widths. For the rest of the eigenvalue computation of the band matrix or after the tridiagonalization, the current best choice is to take advantage of the high parallelism of Cuppen's divide-and-conquer method. Since the implementation of the divide-and-conquer method originates from mapping tasks into binary trees, data mobility is desirable, and the respective tree structures match the actual network contiguity in a distributed parallel configuration. On the other hand, dynamic load balancing on the part of the eigenvalues, which varies depending on the degree of overlap, is important. In this report, we analyze the performance of the divide-and-conquer method in the current implementation of EigenExa and demonstrate the possibility of further performance improvement with respect to speed, scalability, accuracy, etc., observing some automatic tuning results.

Toshiyuki Imamura
RIKEN Center for Computational Science
imamura.toshiyuki@riken.jp

Yusuke Hirota, Mikihiro Hayashi
University of Fukui
y-hirota@u-fukui.ac.jp, hb200839@g.u-fukui.ac.jp

## MS53

### Auto-Tuning for Quantum Computing Related Technology on Supercomputers

To establish high performance computing, performance tuning is still high cost in viewpoint of software productivity. Autotuning (AT) technology is one of promising ways to solve the productivity problem. In this presentation, a state-in-the-art AT technology is presented in several cases for adaptation of AT. In particular, quantum computing related technology is one of hot-topics for next generation computing, but there are many performance parameters to be tuned. We are focusing on adaptation of AT technology to the new target. The adaptation of the AT includes related process for quantum computing on GPU supercomputer environments, such as inspired quantum annealer and simulation of quantum circuits. Preliminary effects of applying AT will be presented.

Takahiro Katagiri
Information Technology Center, Nagoya University
RIKEN R-CCS
katagiri@cc.nagoya-u.ac.jp

## MS53

### Experimenting with the Automatic Tuning of Numerical Software

The use of software libraries on high performance computing systems often requires setting parameters that may greatly impact the performance of the libraries, and consequently of scientific and engineering applications built upon the libraries. In most cases, determining optimal values for the relevant parameters in an impromptu way is impractical. Therefore, the automatic tuning autotuning of such parameters is an activity of great interest. Various tools exist for software autotuning. In this presentation, we summarize the state-of-the-art, and report on lessons learned while using GPTune, an autotuning framework developed by DOEs Exascale Computing Project, for improving the performance of a set of libraries of interest.

Osni A. Marques
Lawrence Berkeley National Laboratory
Berkeley, CA
oamarques@lbl.gov

Makoto Morishita
Graduate School of Informatics,
Nagoya University
morishita@hpc.itc.nagoya-u.ac.jp

## MS53

### Data Analysis Toward Identification of Organizational Models for Software Ecosystem Sustainability

In this presentation the authors describe findings from data analysis methods used to identify organizational models as articulated by the activities of open source software project development and sustainment. In particular, the exascale computing project (ECP) is presented as a case study for how organizational infrastructure can serve software sustainability through key contributions focused on software ecosystem sustainability.

Elaine M. Raybourn
Sandia National Laboratories
emraybo@sandia.gov

Shuji Morisaki
Graduate School of Informatics,
Nagoya University
morisaki.shuji.k2@f.mail.nagoya-u.ac.jp

Killian Muollo
Sandia National Laboratories
kmuollo@sandia.gov

## MS54

### How Mixed Precision Can Accelerate Sparse Solvers

Within the past years, hardware vendors have started designing low precision special function units in response to the demand of the Machine Learning community and their demand for high compute power in low precision formats. This motivates significant research on the design and implementation of mixed precision algorithms. In this talk, we provide an overview of mixed precision numerics and highlight strategies that have proven to be successful in accelerating algorithms without impacting the numerical

properties.

Hartwig Anzt
University of Tennessee, U.S.
hanzt@icl.utk.edu

Terry Cojean, Pratik Nayak, Thomas Gruetzmacher,
Yu-Hsiang Mike Tsai, Marcel Koch, Tobias Ribizel, Fritz
Goebel
Karlsruhe Institute of Technology
terry.cojean@kit.edu, pratik.nayak@kit.edu,
thomas.gruetzmacher@kit.edu, yu-hsiang.tsai@kit.edu,
marcel.koch@kit.edu, tobias.ribizel@kit.edu,
fritz.goebel@kit.edu

Gregor Olenik
Steinbuch Centre for Computing (SCC)
Karlsruher Institut für Technologie (KIT)
gregor.olenik@kit.edu

## MS54

### Pipelined Sparse Solvers: Can More Reliable Computations Help Us to Converge Faster?

Parallel implementations of Krylov subspace methods often help to accelerate the procedure of finding an approximate solution of a linear system. However, such parallelization coupled with asynchronous and out-of-order execution often enlarge the non-associativity impact in floating-point operations. These problems are even amplified when communication-hiding pipelined algorithms are used in order to improve the scalability of Krylov subspace methods. We propose a general framework for deriving robust (reproducible and accurate) variants of Krylov subspace methods. The proposed algorithmic strategies are reinforced by programmability suggestions to assure deterministic and accurate executions. In this talk, we show if the residual replacement strategy in the pipelined Krylov subspace solvers can be mitigated by the higher precision. We illustrate the framework on a few such solvers aiming to restore floating-point associativity and to cope with the attainable precision loss in pipelined solvers.

Roman Iakymchuk
Umeå University
riakymch@cs.umu.se

## MS54

### Residual Inverse Formulation of the Feast Eigenvalue Algorithm Using Mixed-Precision and Inexact System Solves

The standard FEAST eigenvalue algorithm is a method for solving the generalized linear eigenvalue problem $Ax = \lambda Bx$. The key operation in the FEAST algorithm is the numerical evaluation of the contour integral:

$$Q = \frac{1}{2\pi i} \oint_{\mathcal{C}} (zB - A)^{-1} B\tilde{X} dz,$$

where $\mathcal{C}$ is a user-defined closed contour in the complex plane that encloses the wanted eigenvalues, and $\tilde{X}$ is a (usually random) initial guess for the vectors spanning the eigenvector subspace. An alternative form for the contour integral was recently proposed, allowing FEAST to be used for solving nonlinear eigenvalue problems. Rather than the integral above, one instead uses:

$$Q = \frac{1}{2\pi i} \oint_{\mathcal{C}} \left( \tilde{X} - T(z)^{-1} R_E \right) (zI - \tilde{\Lambda})^{-1} dz,$$

where $\tilde{\Lambda}$ is a diagonal matrix of Ritz values corresponding to the approximate eigenvectors in the columns of $\tilde{X}$, $T(z)$ is the residual operator for the eigenvalue problem, and $R_E$ is a matrix whose column vectors are the eigenvector residuals for the approximate eigenpairs $\tilde{X}$ and $\tilde{\Lambda}$. Not only does this residual inverse variant of the FEAST integral find application in addressing non-linear eigenvalue problems, but we will also illustrate its substantial benefits in solving linear eigenvalue problems with mixed-precision and inexact system solves.

Ivan Williams
University of Massachusetts, Amherst
inwilliams@umass.edu

Eric Polizzi
University of Massachusetts, Amherst, USA
epolizzi@engin.umass.edu

## MS54

### Algebraic Programming for High Performance Auto-Parallelised Solvers

Algebraic Programming, or ALP for short, requires programmers to annotate their programs with explicit algebraic information. This information is then used in auto-parallelisation and other automatically applied optimisations, ranging from low-level concerns such as vectorisation to more complex algorithmic transformations. GraphBLAS is a recent framework that allows such algebraic programming, with, as the name implies, a primary focus on graph computing. In this talk, we focus on re-purposing recent progress by ALP/GraphBLAS for solvers. For example, in recent work we demonstrate significant performance improvements due to non-blocking execution, where linear algebraic primitives are lazily evaluated and fused at run-time, speeding up computations by 2.4x and 7.0x versus the comparable frameworks of SuiteSparse:GraphBLAS and Eigen, respectively. These advances naturally benefit iterative solvers as well. After introducing the core concepts and ideas of algebraic programming, the talk will expand on the applicability and performance results of sparse iterative solvers. It then follows with exploring classical direct solvers expressed using ALP, extending the interface with dense linear algebra, structures and views. We sketch our dispatching mechanism with optional JIT-like capabilities for fused execution, and demonstrate some preliminary performance results.

Albert Jan N. Yzelman
Huawei Zürich Research Center
albertjan.yzelman@huawei.com

Denis Jelovina
Computing Systems Lab
Huawei Zürich Research Center
denis.jelovina@huawei.com

Aristeidis Mastoras
Zurich Research Center, Huawei Technologies Switzerland
Computing Systems Laboratory
aristeidis.mastoras@huawei.com

Alberto Scolari, Daniele Giuseppe Spampinato
Computing Systems Lab
Huawei Zürich Research Center
alberto.scolari@huawei.com,
daniele.giuseppe.spampinato@huawei.com

## MS55

### Modern Scientific Workflows in the Era of an Integrated Research Infrastructure

As science marches forward, its dependence on complex computational processes becomes more pronounced. Central to these processes are scientific workflows, serving as critical orchestrators for expansive experiments ranging from cloud-based data preprocessing to intricate multi-facility computational frameworks. As the computing continuum evolves, addressing the needs of contemporary scientific applications necessitates a recalibration of existing systems and the innovation of new workflow functionalities. This presentation offers a panoramic view of cutting-edge advancements in the realm, encapsulating both the general landscape and his own pivotal research contributions. The discussion culminates in pinpointing pressing challenges awaiting solutions within the workflows community.

Rafael Ferreira da Silva
Oak Ridge National Laboratory
silvarf@ornl.gov

## MS55

### End to end Performance Analysis of Complex Workflows

High-Performance Computing (HPC) applications are evolving from the traditional bulk-synchronous parallel paradigms to complex, coordinated workflows with multiple distinct components that leverage AI/ML. Such workflows often have multiple dependencies and distinct binaries, making it challenging to understand their critical path and their performance bottlenecks. Exploring the large configuration space of such workflows for fine-grained performance tuning is equally challenging. In this talk, we will present PerfFlowAspect, an aspect-oriented, low-overhead and open-source tool designed to analyze the performance of complex workflows. We will present a detailed analysis of the applicability of this tool using a modern workflow such as the Autonomous Multi-Scale (AMS) workflow at Lawrence Livermore National Laboratory.

Tapasya Patki
Lawrence Livermore National Laboratory
patki1@llnl.gov

## MS55

### Provenance and Trust in Scientific Workflows

Workflows or pipelines guide and pervade science because they reduce the burden of manual and repeated execution of sometimes trivial steps. Too, as the data and computational needs of models and analysis techniques grow, there is an increasing reliance on larger infrastructure in order to carry out analysis. Workflows can reduce the burden of transitioning from small scale computational environment to large scale environments. Virtual teams working in the academic setting who carry out research that contributes to national cyberinfrastructure face a set of competing challenges when transparency in the infrastructure is critical for determinations of trust of results. For the research team that is exploring AI innovation in infrastructure, these challenges come into high relief when examined through the lens of responsibility to society. What form does responsibility for downstream uses of AI take? We identify the competing interests, and show how they raise challenges when analyzed through the lens of Democratizing AI. Our background is our now 2-year involvement in the NSF AI Institute Intelligent Cyberinfrastructure with Computational Learning in the Environment (ICICLE)

Beth Plale
Indiana University
plale@indiana.edu

## MS55

### Get What You Want/WANDS - Wide Area Network Data Steaming

Large scientific experiments that have a need for data analysis require the compute power available at modern supercomputers. Data generation and computing site are typically not co-located and transporting the data between these facilities is currently a bottleneck. This becomes even more exacerbated when ad-hoc analysis of experiment data is essential. One such use-case is nuclear fusion. For fusion to be a viable concept for energy generation, it is imperative that reactors can respond quickly and correctly to chaotic plasma events to avoid losing confinement. This requires analyzing data gathered at previous shots. To keep the intershot analysis time as short as possible the compute power of modern supercomputers is necessary, creating the need to transfer data across Wide Area Networks efficiently. This talk will present a Wide Area Network Data Streaming (WANDS) framework currently under development. WANDS provides a flexible, user-friendly, and efficient implementation of a server-client model to request subsets of data for processing. The WANDS client has a Python interface for direct integration into data analysis scripts. A lightweight C++ server component on the data production site serves these requests. Communication is facilitated using ADIOS2. A client-side cache further accelerates subsequent requests for previously transferred data.

Stefanie Reuter
Cambridge Open Zettascale Lab
sr2003@cam.ac.uk

## MS56

### Scaling of IPPL: a Massively Parallel, Performance

### Portable C++ Library for Particle-Mesh Methods

We present the Independent Parallel Particle Layer (IPPL), a performance portable C++ library for particle-in-cell methods. IPPL makes use of Kokkos (a performance portability abstraction layer), HeFFTe (a library for large scale FFTs), and MPI (Message Passing Interface) to deliver a portable, massively parallel toolkit for particle-mesh methods. IPPL supports simulations in one to six dimensions, mixed precision, and asynchronous execution in different execution spaces (e.g. CPUs and GPUs). We showcase the performance of the latest version of IPPL using examples from charged particle dynamics (https:// arxiv.org/abs/2205.11052) on state-of-the-art high-performance computing resources, such as Perlmutter and Frontier.

Sonali Mayani, <u>Andreas Adelmann</u>
Paul Scherrer Institut
sonali.mayani@psi.ch, andreas.adelmann@psi.ch

Antoine Cerfon
Courant Institute NYU
cerfon@cims.nyu.edu

Matthias Frey
University of St Andrews
mf248@st-andrews.ac.uk

Veronica Montanaro
ETH
vmontanaro@student.ethz.ch

Sriramkrishnan Muralikrishnan
Juelich Supercomputing Centre
Forschungszentrum Juelich GmbH
s.muralikrishnan@fz-juelich.de

Alessandro Vinciguerra
ETH Zurich
vincigua@student.ethz.ch

### MS56

### Global Impurity Transport for Fusion Devices on Frontier

We describe a new, performance portable PIC-like implementation of global impurity transport (GITR) for magnetically confined fusion devices built on the Cabana and Kokkos libraries. Using a trace impurity assumption, significant numbers of particles and even device-scale simulations to investigate erosion, ionization, and migration are possible. Performance and scalability for the particle integration, background collisions, self-collisions, and redistribution (communication and removal), as well as grid communication and grid-particle interpolation for output distributions are discussed. Initial work with unstructured embedded surfaces for complex systems (and the performance implications) is also described. Performance on both CPU and GPU, including the Frontier system at ORNL, will be shown. Finally, extensions and updates to the Cabana library driven by this work will be discussed.

<u>Lance Bullerwell</u>, Sam Reeve, Kwitae Chong
Oak Ridge National Laboratory
bullerwellle@ornl.gov, reevest@ornl.gov, chongk@ornl.gov

Wenjun Ge
Oak Ridge National Laboratory, U.S.
gew1@ornl.gov

Timothy Younkin, Austin Isner
Oak Ridge National Laboratory
younkintr@ornl.gov, isnerab@ornl.gov

Stuart Slattery
Computational Sciences and Engineering Division
Oak Ridge National Laboratory
slatterysr@ornl.gov

### MS56

### Efficiency of a Parareal Algorithm for Highly Oscillatory Systems Arising in Plasma Physics

We develop a specific version of the parareal algorithm for solving multiscale in time Vlasov-Poisson systems. We use reduced models, obtained from the two-scale convergence theory, for the coarse solving. These models are useful to approximate the original Vlasov-Poisson model at a low computational cost since they are free of high oscillations. The equations are numerically solved with a particle-in-cell method. We provide a thorough analysis of the parallel capabilities of the algorithm on the basis of the ideal speedup. Applications to realistic problems in plasma physics illustrate the efficiency of the proposed approach.

<u>Sever Hirstoaga</u>
Inria Paris
sever.hirstoaga@inria.fr

Pierri-Henri Tournier
LJLL Sorbonne Universite
pierre-henri.tournier@sorbonne-universite.fr

### MS56

### A Deterministic Verification Strategy for Gyrokinetic Particle-in-Cell Algorithms

As particle-in-cell (PIC) algorithms for plasma simulation continue to increase in scope and complexity, a rigorous and straightforward method for verifying PIC implementations is desirable to ensure their correctness. In this presentation, we introduce a deterministic method for the rigorous verification of multidimensional, electrostatic, gyrokinetic particle-in-cell codes based on the method of manufactured solutions. We show that rigorous verification is possible through the exclusive examination of errors in grid quantities, allowing for a very light-weight and non-intrusive implementation in existing particle-in-cell codes. We present numerical results run in XGC that confirm our theoretical claims.

<u>Paul Tranquilli</u>
Lawrence Livermore National Lab
tranquilli1@llnl.gov

Lee Ricketson

Lawrence Livermore National Laboratory
ricketson1@llnl.gov

Benjamin Sturdevant
Princeton Plasma Physics Laboratory
bsturdev@pppl.gov

Luis Chacon
Los Alamos National Laboratory
chacon@lanl.gov

## MS57

### Complexity Reduction of Density Functional Theory Eigenvalue Problems

In Kohn-Sham Density Functional Theory (DFT) calculations, generating the single-particle density matrix is a key computational bottleneck. Conventionally, for molecular systems this has been achieved by diagonalizing the Hamiltonian matrix using dense solvers. This method, however, has a computational cost that scales cubically with the system size, making it less practical for large systems. To counteract this, "diagonalization free" methods have been developed that have costs that scale only linearly with the system size. While these methods improve computational efficiency, they lose details about the Kohn-Sham orbitals and orbital energies, which are important for analyzing systems. To address this limitation, we introduce two new algorithms that utilize information from the density matrix to reduce the cost of computing eigenvalues and eigenvectors of the Hamiltonian. We demonstrate that these algorithms can efficiently be applied to systems composed of thousands of atoms and as an analysis tool that reveals the origin of spectral properties.

William Dawson
Riken Advanced Institute for Computational Science
william.dawson@riken.jp

Eisuke Kawashima
RIKEN Center for Computational Science
eisuke.kawashima@riken.jp

Luigi Genovese
Institut Nanosciences et Cryogénie, CEA
luigi.genovese@cea.fr

Takahito Nakajima
RIKEN Center for Computational Science
nakajima@riken.jp

## MS57

### Novel Algorithms and Hot Electrons to Break the Exascale Barrier

We present the submatrix method and a novel linear-scaling electronic-structure method in conjunction with approximate computing, as well as the implementation of the technique in CP2K. Even though initially proposed for inverse p-th roots, it has recently been recognized that the submatrix method represents a general method to approximate arbitrary matrix functions such as the matrix-sign function of large sparse matrices. The Matrix-sign function is the essential workhorse of linear-scaling electronic-structure theory, and we present an intuitive chemical justification for the accuracy of the submatrix method. We will discuss the efficient implementation of the submatrix method into CP2K with a special focus on limiting communication between compute nodes. The resulting compute kernel is the sign function of a relatively small but dense matrix. Our optimized implementation with a simple diagonalization-based evaluation of the sign function of the submatrices outperforms the Newton-Schulz Sign iteration in initial results, especially for larger cutoffs of matrix elements. This observation shows that the submatrix method will be a valuable tool in the context of approximate computing.

Thomas D. Kühne
Chair of Theoretical Chemistry
Paderborn Center for Parallel Computing
t.kuehne@hzdr.de

## MS57

### Scalability Study for Planewave DFT Solvers

The computer architecture landscape is constantly changing. For example, most modern supercomputers have shifted from the classical computing capabilities offered by CPUs and have adopted the more powerful GPUs, i.e. Summit and Frontier machines are some of the most modern supercomputers that use GPUs to accelerate computation. Given this trend, scientific software needs to adapt to the new features and capabilities offered by these new systems. Therefore, in this talk we focus on planewave Density Functional Theory application, an DOE critical application. More specifically, we present a new distributed Fourier transform implementation that works both on CPUs and GPUs. We outline the capabilities of this new library and present results in the context of the planewave DFT code Quantum Espresso. We integrate our new Fourier transform implementation within the software package and investigate the scalability of the most common algorithms like Conjugate Gradient, PPCG, Davidson and Unconstrained algorithms. We show results for all algorithms targeting Summit and Frontier systems and analyze the scalability of each algorithm providing insights into the main bottlenecks of the computation.

Doru Thom Popovici, Mauro Del Ben, Andrew M. Canning
Lawrence Berkeley National Laboratory
dtpopovici@lbl.gov, mdelben@lbl.gov, acanning@lbl.gov

Osni A. Marques
Lawrence Berkeley National Laboratory
Berkeley, CA
oamarques@lbl.gov

## MS57

### Unitary Canonicalization of Subspace Bases

We demonstrate how to determine a canonical basis for a subspace of a (complex) vector space equipped with a fixed basis. The canonical basis is defined as the basis whose expansion in the host basis is maximally-sparse; the

particular definition used here maximizes the $\ell^4$ norm of the expansion coefficients. Standard nonlinear optimization techniques can then be used to determine the canonical basis. The approach is useful for defining a canonical basis in a variety of generic applications of linear algebra (e.g., for producing unique sets of eigen/singular vectors in presence of degeneracies). Its utility in the context of physical simulation will be illustrated for producing canonical orbitals in systems with discrete and continuous nonabelian symmetries.

Edward Valeev
Virginia Tech
Department of chemistry
efv@vt.edu

## MS58

### Space-Time Multigrid Methods for Convection-Diffusion Equations Arising from Flow Problems

Usually time dependent evolution equations are solved time step by time step where in each step a system of equations corresponding to the spatial discretization has to be solved. Such methods can only be parallelized in space, but the size of the spatial problem limits the strong scaling behavior. Using multigrid methods that treat multiple time steps in an all-at-once system the parallel scaling can be improved significantly in comparison to geometric multigrid solvers in a time stepping application. Due to the improved communication pattern between the parallel processes, this holds true even if a time-simultaneous multigrid method without temporal parallelization is used [Dünnebacke et al., 2021]. Here, we numerically analyze how multigrid waveform relaxation methods [Lubich and Ostermann, 1987] and space-time multigrid methods [Gander and Neumüller, 2014] behave for convection-diffusion(-reaction) equations found in flow problems, when the problems are increasingly transport dominated. The emerging difficulties in such cases can be mitigated by suitable stabilization techniques. Furthermore, we consider problems with space and time dependent diffusion and convection parameters, that arise in global-in-time solution strategies of the Navier-Stokes equations with non-Newtonian fluids.

Jonas Duennebacke
TU Dortmund
Department of Mathematics
jonas.duennebacke@math.tu-dortmund.de

Stefan Turek
TU Dortmund University
Department of Mathematics
stefan.turek.mathematik.tu-dortmund.de

## MS58

### Augmented Lagrangian Acceleration of Global-In-Time Oseen Solvers

This talk focuses on an accelerated global-in-time Oseen solver, which highly exploits the augmented Lagrangian approach to improve the convergence behavior of the Schur complement iteration. The main idea of the solution strat-

egy is to block the individual linear systems of equations at each time step into a single saddle point problem. By elimination of all velocity unknowns, the resulting pressure Schur complement (PSC) equation can be solved efficiently on modern hardware architectures using a space-time multigrid algorithm [Lohmann and Turek, 2023]. However, the accuracy of the involved PSC preconditioners deteriorates as the Reynolds number increases and, hence, causes convergence issues. To improve the robustness of the solution strategy and accelerate its convergence behavior, the augmented Lagrangian approach is exploited by modifying the velocity system matrix in a strongly consistent manner. While the introduced discrete grad-div stabilization does not modify the solution, the accuracy of the adapted PSC preconditioners drastically improves and, hence, guarantees a rapid convergence. This strategy comes at the cost that the involved auxiliary problem for the velocity field becomes ill conditioned so that standard iterative solution strategies are no longer efficient. This calls for highly specialized multigrid solvers, which are based on modified intergrid transfer operators and block diagonal preconditioners (cf. [Benzi and Olshanskii, 2006; Wechsung, 2019]).

Christoph Lohmann
Dortmund University of Technology
christoph.lohmann@mathematik.tu-dortmund.de

Stefan Turek
TU Dortmund University
Department of Mathematics
stefan.turek.mathematik.tu-dortmund.de

## MS58

### A Parallel-in-Time Spectral Deferred Correction Finite Element Method for Unsteady Incompressible Viscous Flow Problems

Simulating unsteady viscous flows by numerically solving the time-dependent Navier-Stokes equations is a computationally expensive challenge. However, while spatial parallelization can reduce computational costs, temporal integration of time-sensitive applications often requires a very large number of time steps. Therefore, more parallelism in numerical time-stepping schemes for further speedup is required. The present work proposes and analyzes a parallel-in-time spectral deferred correction method for the solution of the unsteady incompressible viscous flow problems governed by parabolicelliptic PDEs. The temporal discretization employs the Spectral Deferred Correction (SDC) method in parallel, which iteratively computes a higher-order collocation solution by conducting a sequence of correction sweeps through the utilization of a low-order time-stepping technique. A standard finite element method is considered for spatial discretization due to its ability to accurately capture complex geometries and boundary conditions. The goal of this work is to illustrate and analyze the properties of the parallel-in-time method through numerical experiments including flows past a cylinder (using the standard DFG 2D-3 benchmark) which is selected as an unsteady flow example.

Abdelouahed Ouardghi
Forschungszentrum Juelich

Juelich Supercomputing Centre
a.ouardghi@fz-juelich.de

Robert Speck
Juelich Supercomputing Centre
Forschungszentrum Juelich GmbH
r.speck@fz-juelich.de

## MS59

### Superchips for Supercomputing: Grace Hopper and Its Impact on Science Use Cases

Grace Hopper Superchips bring a new level of integration to supercomputing. In this talk I will detail these new capabilities, explain how they impact applications and enable user productivity. With the end of Moores law innovations in power efficiency, the ability to accelerate more applications and performance improvements found in Grace Hopper and its lowering of the programming bar are key to further advancing application performance and reducing data center energy usage.

Ian Karlin
NVIDIA
ikarlin@nvidia.com

## MS59

### Improving and Extending Collectives and Rma for New Architectures

New Accelerated Processing Unit have emerged, increasing further the memory hierarchy of a HPC node. Meanwhile, new form of HPC system architecture have been devised. An example of such architecture is called Modular system architecture, in which multiple modules of different architecture and capabilities will be interconnected together, to better adapt to the different types of computation. The resulting network topology increases in complexity. To take advantage of the performance of these new architectures, extension and adaptation of MPI runtime are required. A first modification is to adapt collectives communications to the HPC hardware and network topology. In this presentation we will exhibit how we extended the existing hierarchical collective capabilities of Open MPI to take into account more topological levels. We will also present our recent effort in improving MPI RMA capabilities. Applications can make use of such communication type as they enable better computation/communication overlap, even if using them is notoriously difficult as relaxing synchronization can lead to memory accesses race conditions.

Pierre Lemarinier
EVIDEN
pierre.lemarinier@eviden.com

## MS59

### Architecting An Apu: Hardware and Software Co-Design for Exascale Systems

The Exascale barrier was first breached with the Frontier supercomputer, which is built on a heterogeneous CPU and GPU architecture. While this system design has enabled several remarkable achievement in high-performance computing, applications run at exascale have already identified multiple opportunities where this paradigm can be improved; notably, the communication costs, and the complexity of the resultant programming model, incurred by the presence of two isolated memory spaces for CPU and GPU. In this talk, we discuss a tight hardware and software co-design effort between AMD, HPE, and Lawrence Livermore National Laboratory, dating back to the DOE Fastforward, Pathforward, and CORAL2 programs, and resulting in the AMD Instinct MI300 APU (Accelerated Processing Unit) architecture. The talk will discuss demonstrated advantages, and future possibilities, afforded by MI300 for Exascale computing, including: the improved simplicity of porting from CPU codes, performance benefits resulting from close integration of CPU and GPU subsystems; simplifications and improvements are realized in a variety of tools, including RAJA and Kokkos accelerator abstraction frameworks, a recently developed Standard Parallelism interface to AMD APUs, and automatic offload systems.

Malaya Nicholas, Michael Rowan
AMD
nicholas.malaya@amd.com, michael.rowan@amd.com

## MS59

### Co-Design Serendipity: Leveraging the Cerebras CS-2 AI Accelerator for Seismic Processing

The recent explosion in the AI market has led to the development of specialized hardware for performing training and inference on neural networks. The Cerebras CS-2 is one such hardware innovation: the system is built from the ground up to accelerate large AI training workloads. The CS-2 is built around the wafer-scale engine (WSE), a single chip consisting of 850,000 processing elements (PEs) laid out in a 2D mesh, with each PE containing 48 kB of memory accessible in a single cycle. While this chip has explicitly been the product of co-design with AI workloads, many of its features also address bottlenecks present in more traditional HPC workloads. We discuss the serendipitous consequences of this co-design process by presenting HPC problems which are massively accelerated by the CS-2. In particular, we present the mapping of a key kernel arising in many wave-equation-based algorithms for seismic processing. This kernel, TLR-MVM, takes advantage of the sparsity in seismic data to reduce the most expensive component of multi-dimensional convolutions (MDC) to a series of embarrassingly parallel matrix-vector multiplications. We execute TLR-MVM for a standard seismic dataset across 48 CS-2s, and achieve sustained memory bandwidth of 92.58 PB/s. Potential opportunities for future HPC workloads on Cerebras systems are also presented.

Leighton Wilson
Cerebras
leighton.wilson@cerebras.net

## MS60

### Analysis of Randomized Householder-Cholesky Qr Factorization with Multisketching

CholeskyQR2 and shifted CholeskyQR3 are two state-of-

the-art algorithms for computing tall-and-skinny QR factorizations since they attain high performance on current computer architectures. However, to guarantee stability, for some applications, CholeskyQR2 faces a prohibitive restriction on the condition number of the underlying matrix to factorize. Shifted CholeskyQR3 is stable but has 50% more computational and communication costs than CholeskyQR2. In this talk, a randomized QR algorithm called Randomized Householder-Cholesky (`rand_cholQR`) is analyzed. Using one or two random sketch matrices, it is proved that with high probability, its orthogonality error is bounded by a constant of the order of unit roundoff for any numerically full-rank matrix, and hence it is as stable as shifted CholeskyQR3. An evluation of the performance of `rand_cholQR` on a NVIDIA A100 GPU demonstrates that for tall-and-skinny matrices, `rand_cholQR` with multiple sketch matrices is nearly as fast as, or in some cases faster than, CholeskyQR2. Hence, compared to CholeskyQR2, `rand_cholQR` is more stable with almost no extra computational or memory cost, and therefore a superior algorithm both in theory and practice.

Andrew J. Higgins
Temple University, U.S.
andrew.higgins@temple.edu

Andrew J. Higgins
Temple University
andrew.higgins@temple.edu

Daniel B. B. Szyld
Temple University
Department of Mathematics
szyld@temple.edu

Erik G. Boman
Center for Computing Research
Sandia National Labs
egboman@sandia.gov

Ichitaro Yamazaki
Sandia National Laboratories
iyamaza@sandia.gov

## MS60

### Randomized Matrix Algorithms Involving Orthogonalization

We investigate the potential of randomization as a viable alternative in the context of orthogonalization, including the computation of orthonormal bases, matrix decompositions, low-rank approximations, and the solution of linear systems and of least squares/regression problems. We consider various sampling and sketching strategies; evaluate matrix concentration inequalities and structural bounds for estimating the error due to randomization; examine numerical stability issues; present empirical results; and identify problem classes where randomization delivers sufficient accuracy and speedup.

Ilse C. F. Ipsen
North Carolina State University
Department of Mathematics
ipsen@ncsu.edu

## MS60

### Randomized Householder QR Factorization with Applications to Krylov Subspace Solvers

The Randomized Householder QR factorization (RHQR) can be used to obtain a well conditioned basis of a set of high dimensional vectors. We discuss its performance and numerical stability, as well as its usage in Arnoldi process and GMRES to solve systems of linear equations. We detail numerical experiments in which both in simple and double precision, RHQR outputs a better conditioned basis than Randomized Gram-Schmidt (RGS) and produces accurate factorizations with very small sampling size.

Edouard Timsit
INRIA
edouard.timsit@inria.fr

Laura Grigori
EPFL and PSI, Switzerland
laura.grigori@epfl.ch

## MS61

### Scalable Algorithms for Particle-In-Cell Plasma Modeling with WarpX

WarpX is a massively parallel electromagnetic and electrostatic particle-in-cell code developed under the US DOE Exascale Computing Project. WarpX was used in the 2022 ACM Gordon Bell Prize for mesh-refined simulations performed on the worlds fastest supercomputers, including the world's first Exascale machine, Frontier (OLCF). WarpX has been used extensively for modeling a variety of physics applications, including laser-plasma interactions, accelerator and beam physics, astrophysical plasma, fusion devices, and microelectronics. In this talk, I will present the scalable algorithms and performance portable numerical methods implemented in WarpX, including pseudo-spectral analytical time-domain method (PSATD), mesh-refinement, dynamic load balancing, and particle sorting modules that play an important role in efficient modeling of kinetic processes in plasma applications on heterogeneous architectures.

Revathi Jambunathan, Arianna Formenti, Marco Garten
Lawrence Berkeley National Laboratory
rjambunathan@lbl.gov, ariannaformenti@lbl.gov,
mgarten@lbl.gov

David Grote
Lawrence Livermore National Laboratory
dpgrote@lbl.gov

Axel Huebl
Lawrence Berkeley National Laboratory (LBNL), USA
axelhuebl@lbl.gov

Hannah Klion, Prabhat Kumar
Lawrence Berkeley National Laboratory
klion@lbl.gov, prabhatkumar@lbl.gov

Rémi Lehe
Lawrence Berkeley National Laboratory
ATAP
rlehe@lbl.gov

Andrew Myers
Lawrence Berkeley National Lab
atmyers@lbl.gov

Ryan Sandberg
Lawrence Berkeley National Laboratory
rsandberg@lbl.gov

Olga Shapoval
LBNL
OShapoval@lbl.gov

Jean-Luc Vay
Lawrence Berkeley National Laboratory
Berkeley, CA
JLVay@lbl.gov

Weiqun Zhang
Lawrence Berkeley National Laboratory
Center for Computational Sciences and Engineering
weiqunzhang@lbl.gov

Edoardo Zoni
Lawrence Berkeley National Laboratory
ezoni@lbl.gov

## MS61
### PIC-DSMC Techniques for Complex, Reacting Multiscale Flows

Simple DSMC simulations are inherently scalable, but as problems become complex and include spatially and temporally varying features such as orders of magnitude differences in number densities or rare species with important chemistry and local nonequilibrium, the cost of simulating small regions can grow enormously and quickly. Various methods have been developed to control local costs, but these can have non-obvious drawbacks that are surprising even to experienced users. Using EMPIRE, Sandias PIC-DSMC code with variable-weight particles for capturing rare species and merge algorithms for controlling particle counts, we demonstrate some unappreciated failure modes including Gamblers Ruin for random merging and weight coalescence for M-to-1 and M-to-2 merges. We emphasize the insufficiency of merely demonstrating convergence in the limit of many particles and small elements while running with relatively few particles and time- or ensemble-averaging.

William McDoniel, Christopher Moore
Sandia National Laboratories
wmcdoni@sandia.gov, chmoore@sandia.gov

## MS61
### DPG-Based Vlasov Solvers with Adaptive Mesh Refinement

Efficient solution of the Vlasov equation, which can be up to six-dimensional, is key to the simulation of many difficult problems in plasma physics. The discontinuous Petrov-Galerkin (DPG) finite element methodology provides a framework for the development of stable (in the sense of LBB conditions) finite element formulations, with built-in mechanisms for adaptivity. We present two DPG-based formulations for Vlasov: a time-marching, backward-Euler formulation, and a space-time formulation, with an ultimate target of solving problems in the full seven-dimensional setting. For this purpose, we employ tensor-product data representations supported by recent additions to the Intrepid2 package within Trilinos, as well as corresponding developments within Camellia, a finite element library designed to facilitate rapid development of computationally efficient, hp-adaptive finite element solvers, starting with support for DPG. In this talk, we discuss our progress to date, including several adaptive mesh refinement results from 1D1V time-marching and space-time Vlasov-Poisson problems.

Nathan V. Roberts, Keith Cartwright, Adam Darr
Sandia National Laboratories
nvrober@sandia.gov, klcartw@sandia.gov, amdarr@sandia.gov

## MS61
### Projection-Based Reduced Order Model for Accelerating Kinetic Simulations of Electrostatic Plasmas

In this talk, we present a computationally efficient method using reduced-order modeling (ROM) for collisionless electrostatic plasma kinetics governed by the Vlasov-Poisson equation. High-performance computing and modern numerical algorithms have made high-fidelity kinetic plasma simulations tractable. Despite continued advances in these areas, Eulerian simulations and particle-based methods are too costly to model fusion reactors, where many configurations and plasma modes must be considered. Our ROM approach projects the equation onto a linear subspace defined by principal proper orthogonal decomposition (POD) modes. We introduce an efficient tensorial method to update the nonlinear term using a precomputed third-order tensor. We capture multiscale behavior with a minimal number of POD modes by decomposing the solution into multiple time windows using a physical-time indicator and creating a temporally-local ROM. Applied to 1D–1V simulations, specifically the benchmark two-stream instability case, our time-windowed reduced-order model (TW–ROM) with the tensorial approach solves the equation approximately 280 times faster than Eulerian simulations while maintaining a maximum relative error of 4% for the training data and 13% for the testing data. This work performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.

Ping-Hsuan Tsai
University of Illinois
pht2@illinois.edu

Seung Whan Chung
Lawrence Livermore National Laboratory
chung28@llnl.gov

Debojyoti Ghosh
Lawrence Livermore National Laboratory
Livermore, CA
ghosh5@llnl.gov

John Loffeld, Youngsoo Choi, Jonathan Belof
Lawrence Livermore National Laboratory
loffeld1@llnl.gov, choi15@llnl.gov, belof1@llnl.gov

## MS62
### Rethinking Heterogeneous Library Design

High performance libraries have traditionally been designed to be run on specific devices. However, as computing architectures evolve to incorporate multiple heterogeneous devices on a single platform, a library now has to be able to seamlessly run on different devices, and potentially on multiple heterogeneous devices. Given the large discrepancy in computing capability, it is often a challenge identifying when a problem should be run on a single or multiple devices. In this talk, we discuss our approach of a common abstraction for CPUs and GPUs and how this abstraction can simplify the design of performance implementations for both CPUs and GPUs.

Tze Meng Low
Carnegie Mellon University
lowt@andrew.cmu.edu

## MS62
### Linear Algebra Software for a Changing Accelerator Landscape

The established split between hardware accelerators and the rest of the modern compute nodes has partially been solved with an ever increasing drive towards tight integration between components. The linear algebra software must keep keep pace with these changes to take advantage of the plethora of floating-point formats but also algorithmic advancements in both mixed-precision and randomized methods. This talk will give an overview of these hardware trends and the software libraries that target them.

Piotr Luszczek
Massachusetts Institute of Technology
luszczek@ll.mit.edu

## MS62
### MatRIS: Multi-Level Math Library Abstraction for Heterogeneity and Performance Portability using IRIS Runtime

Vendor libraries are tuned for one architecture and are not portable to others. Moreover, these lack support for heterogeneity and multi-device computation orchestration which is required for efficiently using contemporary HPC and Cloud resources. To address these challenges, we introduce MatRIS, a multi-level abstraction for scalable and performance portable math library for sparse/dense BLAS/LaPACK operations using IRIS runtime. MatRIS-IRIS co-design introduces three levels of abstraction to make the implementation completely architecture agnostic, providing high programming productiv-

ity. We demonstrate that MatRIS is portable without any change in source code and can fully utilize different multi-device heterogeneous systems, achieving high performance and scalability on three heterogeneous systems, Summit, Frontier, and a CADES cloud node with four NVIDIA A100 GPUs and four AMD MI100 GPUs. A detailed performance study is presented where MatRIS demonstrates multi-device scalability. When compared, MatRIS provides competitive and even better performance with respect to vendors and other libraries.

M.A.H Monil
Oak Ridge National Laboratory, USA
monilm@ornl.gov

Pedro Valero-Lara, Keita Teranishi, Narasinga Rao Miniskar
Oak Ridge National Laboratory
valerolarap@ornl.gov, teranishik@ornl.gov, miniskarnr@ornl.gov

Jeffrey S. Vetter
Oak Ridge National Laboratory
University of Tennessee, Knoxville
vetter@computer.org

## MS62
### Towards a Unified Micro-kernel Abstraction for GPU Linear Algebra

We have created a micro-kernel abstraction for GPUs robust enough to represent the tensor core and data movement operations from NVIDIA GPU architectures spanning Maxwell all the way to Hopper. In this talk, we discuss how CuTe layouts and layout algebra allow us to uniformly represent GPU architecture specific operations in a consistent programming model regardless of the threads and data they operate upon to build CUTLASS 3.xs core abstractions.

Vijay Thakkar
Georgia Tech
NVIDIA
thakkarv@gatech.edu

Rich Vuduc
Georgia Institute of Technology
richie@cc.gatech.edu

## MS63
### A Multi-Model Parallel Algorithm for Explainable Learning

The task of modelling data for interpretation and/or prediction is usually focused on selecting the best single model (learner) that fits the data and that predicts future data accurately. In many cases though this approach has either led to an increased use of black-box machine-learning mechanisms which are hard to interpret for practitioners (and can be prone to overfitting) or to model- and variable-interpretations that vary between analysts even when using the same data. This has resulted in the so-called "replication crisis" where, in different domains, findings from dif-

ferent studies appear to contradict each other. To address this problem and improve interpretation and individual-level prediction, a new line of research fosters the use of many models (so-called "Rashomon Sets"). Developed in parallel to this direction of research, we present a heuristic algorithm which aims at selecting a set of highly predictive models which can deliver new ways of (i) interpreting data, (ii) allowing practitioners to pick the highly-predictive model that best suits their needs, (iii) delivering new forms of inference when there is more than one true model, and (iv) providing a method to eventually improve ensemble learning methods. In particular, we discuss how this new approach can be parallelized for efficient estimation of the Rashomon Set and how it can be insightful in the domain of medicine.

Roberto Molinari
Auburn University
Department of Mathematics and Statistics
robmolinari@auburn.edu

Stephane Guerrier
University of Geneva
stephane.guerrier@unige.ch

Nabil Mili
University of Lausanne
nabil.mili@unil.ch

Cesare Miglioli
University of Geneva
cesare.miglioli@unige.ch

Gaetan Bakalli
EM Lyon Business School
bakalli@em-lyon.com

Yagmur Yavuz Ozdemir
Auburn University
yzy0096@auburn.edu

## MS63

### Robustness in Deep Learning

Deep neural networks are playing increasing roles in machine learning and artificial intelligence. Their performance highly depends on the network architecture and the loss function. The classical square loss is widely known to be sensitive to outliers. We propose the use of robust loss and two stage algorithms for deep neural networks, which are able to extract robust features and deal with outliers effectively. Applications in regression analysis and adversarial machine learning will be discussed.

Shu Liu
Middle Tennessee State University
sl6b@mtmail.mtsu.edu

Qiang Wu
Department of Mathematical Sciences
Middle Tennessee State University
qiang.wu@mtsu.edu

## MS64

### Testing for Floating-Point Exception Errors when Porting to Heterogeneous Systems

As scientific codes are ported to heterogeneous systems, numerical correctness, and reproducibility remain crucial. Incorrect handling floating-point exceptions such as Infinity, not a number (NaN), and subnormal numbers can cause numerical reproducibility issues and impact numerical correctness. We present two techniques to supplement testing numerical software for floating-point exceptions: GPU-FPX, a tool to detect exceptions at runtime in GPUs via binary instrumentation, and Xscope, a framework based on Bayesian optimization to find unknown inputs that trigger exceptions in black box programs (where the source is not available).

Ignacio Laguna
Lawrence Livermore National Laboratory
ilaguna@llnl.gov

Ganesh Gopalakrishnan, Xinyi Li
University of Utah
ganesh@cs.utah.edu, u1266620@utah.edu

## MS64

### OpenMP Validation and Verification Test-suite Insights and Lessons Learnt

The Validation and Verification (VV) effort for OpenMP has been critical for ECP as both applications and implementations are simultaneously developed and enhanced to match the newer and much more powerful Exascale machines. Over the duration of the ECP project we developed tests and extracted OpenMP application kernels for independent verification of OpenMP implementations. Our test-suite ensures that critical OpenMP features are rigorously tested by regularly polling application projects to track most used/desired features and providing then providing extensive coverage using combinations of OpenMP directive and clauses tests. This helped uncover issues and regressions in vendor implementations and helped fix numerous implementation bugs that were not caught by their in-house testing efforts. In this talk we would like to discuss our approach, findings and lessons learnt over the course of this project.

Swaroop Pophale
Oak Ridge National Laboratory
pophaless@ornl.gov

Sunita Chandrasekaran
University of Delaware
schandra@udel.edu

## MS64

### Toward Automatic Test Synthesis for Kokkos Performance Portable Programming Models

Performance Portable Programming Models have significantly improved productivity when it comes to lever-

aging various heterogeneous node architectures in high-performance computing (HPC). Nevertheless, these portability layers still contain intricacies in how they interact with target node architectures and their low-level runtime systems. This often demands expert knowledge from application developers to write programs that are both portable and free of errors. To address this challenge, we propose the use of program analysis tools and methodologies for Kokkos, a prominent performance portable programming framework. This approach eliminates the need for platform-specific testing to ensure the correctness of applications developed using performance portable programming frameworks. Firstly, we will introduce an automated analysis tool designed for Kokkos parallel programs. This tool employs guided symbolic execution through an LLVM-based Klee plugin. This automated analysis serves as an initial pass over the compiled program, aiding in the detection of bugs and acting as a cost-effective precursor to dynamic analysis. Secondly, we will explore how a collection of parallel programming examples, gathered from the community and categorized as either correct or incorrect, can enhance the feasibility of concolic analysis for parallel programs. [SNL is managed and operated by NTESS under DOE NNSA contract DE-NA0003525.]

Keita Teranishi
Oak Ridge National Laboratory
teranishik@ornl.gov

Richard Rutledge
Sandia National Laboratories
rlrutle@sandia.gov

Vivek Kale
Brookhaven National Laboratory
vivek.lkale@gmail.com

Samuel Pollard, Shyamali Mukherjee, Jackson Mayo
Sandia National Laboratories
spolla@sandia.gov,                smukher@sandia.gov,
mayo@sandia.gov

## MS64

### Determinacy Checking of Heterogeneous Parallel Programs

The end of Moore's Law has led to multiple hardware disruptions for future HPC systems, including a trend towards extreme heterogeneity for sustained performance improvements. This talk will focus on the new correctness challenges faced by application developers who aim to use portable programming systems like OpenMP and Kokkos to target a broad range of heterogeneous hardware. These challenges include new forms of data race and memory consistency errors that can led to nondeterministic results when the same program is executed on different heterogeneous processors. In addition to summarizing the challenges, this talk will discuss the state of the art in static and dynamic approaches that aim to automate determinacy checking for heterogeneous parallel programs.

Lechen Yu, Vivek Sarkar
Georgia Institute of Technology
lechen.yu@gatech.edu, vsarkar@gatech.edu

## MS65

### Parallel, Portable Sparse Code Generation with MLIR and Kokkos

The MLIR (multi-level intermediate representation) in LLVM can represent operations ranging from high-level (e.g. matrix multiplication) to low-level (e.g. a bitwise OR). It also includes automatic transformations for lowering high-level operations to equivalent sequences of low-level operations. One active area of development in MLIR is support for sparse tensors as a fundamental data type. In this work, we demonstrate an MLIR-based compiler pipeline that can take a high-level Python program with sparse tensor expressions, and turn it into parallel C++ source code using the Kokkos programming model. This source code can then be incorporated into an existing Kokkos-based application, and compiled for any Kokkos backend (including OpenMP, Cuda and HIP).

Brian Kelley, Kim Liegeois, Siva Rajamanickam
Sandia National Laboratories
bmkelle@sandia.gov,         knliege@sandia.gov,         srajama@sandia.gov

## MS65

### Tackling the Challenges of High-Performance Graph Analytics at Compiler Level

This presentation introduces COMET, a DSL and compiler framework for dense and sparse tensor algebra that employs progressive lowering to generate efficient code for target heterogeneous systems starting from several high-level languages. COMET has been designed to solve some of the challenges in scientific, data analytics, and AI workloads, commonly used at National Laboratories. In addition to common compiler optimizations and code transformations, COMET employs domain-specific and architecture-specific optimizations leveraging the semantics expressed by high-level languages and architectural features. This talk specifically focuses on high-performance data analytics, highlighting challenges and opportunities for domain-specific optimizations at the compiler level and describing the optimizations and code transformation employed by COMET to generate efficient code for graph algorithms.

Gokcen Kestor, Rizwan Ashraf, Zheng Peng, Polykarpos Thomadakis
Pacific Northwest National Laboratory
gokcen.kestor@pnnl.gov,         rizwan.ashraf@pnnl.gov,
zhen.peng@pnnl.gov, polykarpos.thomadakis@pnnl.gov

## MS65

### Challenges in Compiler Optimization of Sparse Matrix/Tensor Computations

Sparsity is ubiquitous in scientific computing applications and is of increasing importance in machine learning. Numerical libraries have been the primary resource for application developers of applications involving sparse matrix/tensor computations. Efficient numerical libraries have traditionally been manually developed by experts.

However, the increasing diversity of hardware platforms and the number of opportunities to improve performance in machine learning pipelines by creating "fused" kernels has necessitated advances in automated compiler optimization of sparse computations. This talk will provide an overview of key challenges in compiler optimization of sparse matrix/tensor computations, to set the stage for the talks in this minisymposium that describe several recent advances on this topic.

Ponnuswamy Sadayappan
University of Utah
saday@cs.utah.edu

## MS65

### Polyhedral Specification and Code Generation of Sparse Tensor Contraction with Co-Iteration

We present a code generator for sparse tensor contraction computations based on a mathematical representation of loop nest computations in the sparse polyhedral framework (SPF). SPF extends the polyhedral model to support non-affine computations, such as arise in sparse tensors. This work extends SPF to perform layout specification, optimization, and code generation of sparse tensor code: 1) we developed a polyhedral layout specification that decouples iteration spaces for layout and computation; and, 2) we developed a code generator that efficiently co-iterate sparse tensors by combining polyhedra scanning using an ILP solver with the synthesis of sparse co-iteration using an SMT solver. This code generator is more versatile to support additional layout and computation combinations. This flexibility is brought by the mathematical abstraction of SPF that allows us to decouple 1) logical layout and physical layout, 2) layout and computation, and 3) computation and execution schedule. This code generator also generates more efficient code when compared to a state-of-the-art tensor compiler, TACO. It achieved, on average, 1.63x faster parallel performance than TACO on sparse-sparse co-iteration. We can even improve that to 2.72x average speedup by switching the co-iteration algorithm generated through the code synthesis.

Tuowen Zhao
University of Utah
ztuowen@gmail.com

Tobi Popoola
Boise State University
tobipopoola@u.boisestate.edu

Mary Hall
School of Computing
University of Utah
mhall@cs.utah.edu

Catherine Olschanowsky
Boise State University
catherineolschan@boisestate.edu

Michelle Strout
University of Arizona

mstrout@cs.arizona.edu

## MS66

### Algorithm Differentiation in Dolfinx

In this presentation, we introduce a new automatic differentiation framework tailored for the DOLFINx environment. This framework has been developed to provide users with enhanced control over their computations, particularly in scenarios where fine-grained control is crucial, such as achieving scalability on high-performance computing (HPC) systems. Given the architectural changes introduced in DOLFINx compared to its predecessor, DOLFIN, we have undertaken a comprehensive redesign of the automatic differentiation framework to seamlessly integrate with the new lower-level API. The new framework leverages custom local computational graphs to establish abstract representations for both the forward and backward models. By relying on wrapping the existing DOLFINx functions and objects, we reduce the obstacles for the user while allowing more experienced users to fully customise and manipulate the computational graph. During our presentation, we will showcase scalability results across diverse problem types, highlighting the framework's adaptability to various computational challenges. Furthermore, we will explore its integration with AI frameworks, such as PyTorch, emphasising the critical significance of this synergy in enabling users to leverage the strengths of both DOLFINx and PyTorch to tackle complex scientific and engineering problems.

Igor Baratta
Department of Engineering
University of Cambridge
ia397@cam.ac.uk

Joergen Dokken
Simula Research Laboratory
dokken@simula.no

Niklas Hornischer
Cambridge University
nh605@cam.ac.uk

## MS66

### Differentiable Programming Across the Pde and Ml Divide

PDEs are central to describing and modelling complex physical systems that arise in many disciplines across science and engineering. However, in many realistic applications, PDE modelling provides an incomplete description of the physics of interest. Machine learning techniques have become increasingly popular in filling the knowledge gap between the fundamental physical laws, expressed as differential equations, and the real-world phenomena studied by scientists and engineers. The emergence of this approach urges the need for scientific simulation frameworks that allow for the efficient development and deployment of models coupling PDEs and ML. We employ a differentiable programming approach to build a highly efficient and composable interface that provides researchers, engineers, and domain experts with diverse backgrounds with

a highly productive way to run high-performance simulations coupling PDEs, implemented using the finite element method (FEM) in Firedrake, and machine learning models, specified in PyTorch. The resulting framework maintains separation of concerns while only requiring trivial changes to existing code.

Nacime Bouziani, Nacime Bouziani
I-X Centre for AI In Science, Imperial College London
n.bouziani18@imperial.ac.uk,
n.bouziani18@imperial.ac.uk

David Ham
Department of Mathematics
Imperial College London
david.ham@imperial.ac.uk

## MS66
### Scientific Machine Learning in JAX

We're building an open-source ecosystem for scientific computing and machine learning in JAX! Differential equation solvers, neural networks, root finding, optimisation, etc. Many scientific problems require computational modelling, and these can now take advantage of the ubiquitous autodifferentiation, autoparallelism, and GPU/TPU acceleration, offered by modern computational frameworks. (Going beyond that offered by older tools like SciPy, MAT-LAB, or Julia.) My talk will offer an overview of this work, and which I hope will provide new tools for you to solve the problems you are tackling.

Patrick Kidger
Google X, Mountain View
kidger@google.com

## MS66
### Differentiable Templates: Composing Simulations with Machine Learning

Ever since the rise of physics-informed neural networks [Raissi et al., 2019], and deep equilibrium models [Bai et al., 2019] of which Neural ODEs [Chen et al., 2018] are the most prominent example, the line between classical parallel scientific simulation, and modern machine learning has been blurred further and further. At times, attempting to embed the entire differentiable solver with the machine learning layers to then end-to-end optimize [Kidger 2022], others integrate the differentiable solver in an outer-loop optimization or analysis procedure [Um et al., 2020]. Not utilizing existing, painstakingly developed ODE-, and PDE-solver infrastructure, these approaches had entirely new differentiable solver ecosystems written for them in PyTorch, and JAX. In this talk we will present our new framework *Differentiable Templates*, which follows a compiler-guided approach to compose between existing verified, high-performance scientific solvers using compiler-based automatic differentiation with Enzyme and PyTorch. Taking production C/C++, and FORTRAN solvers we will showcase the efficacy of our framework by integrating trained machine learning models into our production solvers, online training machine learning modules inside of our solvers, and exporting entire solvers into PyTorch

strictly through our compiler analyses and transformations.

Ludger Paehler, Tim Gymnich
Technical University of Munich
ludger.paehler@tum.de, tim.gymnich@tum.de

Jan Hueckelheim, Shk Narayanan
Argonne National Laboratory
jhueckelheim@anl.gov, snarayan@mcs.anl.gov

Josef Winter
Technical University of Munich
josef.winter@tum.de

William S. Moses
Massachusetts Institute of Technology
wsmoses@illinois.edu

Johannes Doerfert
Lawrence Livermore National Laboratory
doerfert1@llnl.gov

Nikolaus Adams
Technical University of Munich
nikolaus.adams@tum.de

## MS67
### Hacc: Extreme Scaling and Performance Across Diverse Architectures

Supercomputing is dominated by heterogeneous computational architectures, characterized by the necessity to exploit concurrency at multiple scales. The HACC (Hardware/Hybrid Accelerated Cosmology Code) framework exploits this diverse landscape at the largest scales of problem size, obtaining high scalability and sustained performance. Developed to satisfy the science requirements of cosmological surveys, HACC melds particle and grid methods for gravitational interactions, gas dynamics, and a variety of subgrid models for astrophysical processes, using a novel algorithmic structure that flexibly maps across architectures, including CPU/GPU and multi/many-core systems.

Salman Habib
Argonne National Laboratory
habib@anl.gov

Nicholas Frontiere
Argonne National Laboratory and University of Chicago
nfrontiere@anl.gov

Katrin Heitmann, Jeffrey Emberson, Michael Buehlmann, Esteban Rangel, Adrian Pope
Argonne National Laboratory
heitmann@anl.gov, jemberson@anl.gov,
mbuehlmann@anl.gov, erangel@anl.gov, apope@anl.gov

Vitali Morozov
ALCF, Argonne National Laboratory

morozov@anl.gov

## MS67

### A Massive Space-Time Parallel Particle-In-Fourier Framework for Kinetic Plasma Simulations

Particle-In-Fourier (PIF) schemes are attractive for long-time integration of kinetic plasma simulations as they conserve charge, momentum and energy, exhibit a variational structure, do not have aliasing and have excellent stability properties. However, they are typically more expensive than the commonly used Particle-In-Cell (PIC) schemes due to the requirement of non-uniform DFTs or FFTs. In this talk, we propose a Parareal-based parallel-in-time integration method for PIF schemes with PIF as the fine propagator and standard PIC scheme as the coarse propagator towards the goal of performing long-time integration simulations with PIF schemes. The resulting scheme is implemented on the performance portable library IPPL and we numerically investigate the convergence of it with respect to the discretization parameters of PIF and PIC. We present space-time parallel simulations with the proposed scheme for a variety of benchmarks such as the Landau damping, the two-stream instability and the Penning trap on thousands of A100 GPUs on the JUWELS Booster supercomputer and show 3.6-4.3X speedup compared to space-only parallelization.

Sriramkrishnan Muralikrishnan, Robert Speck
Juelich Supercomputing Centre
Forschungszentrum Juelich GmbH
s.muralikrishnan@fz-juelich.de, r.speck@fz-juelich.de

## MS67

### PUMIPic: Parallel Unstructured Mesh Infrastructure for PIC

Anisotropic unstructured mesh discretizations enable effective resolution of key features in geometrically complex plasma physics devices such as ITER, DIII-D, and stellarators. Running multiscale particle-in-cell (PIC) simulations on these meshes requires a scalable, easy to use PIC infrastructure that can be tuned to different system architectures. PUMIPic, the parallel unstructured mesh infrastructure for PIC, is a C++, MPI+Kokkos based, library that meets these needs. An overview of PUMIPic, its methods, and supporting workflow components, will be presented along with progress on developing and verifying the edge plasma physics application, XGCm, and the impurity transport application, GITRm.

Cameron W. Smith, Angel Castillo
Scientific Computation Research Center
Rensselaer Polytechnic Institute
smithc11@rpi.edu, castia5@rpi.edu

Chonglin Zhang
University of North Dakota
chonglin.zhang@und.edu

Dhyanjyoti Nath
Scientific Computation Research Center
Rensselaer Polytechnic Institute

nathd@rpi.edu

Onkar Sahni
Rensselaer Polytechnic Institute
sahni@rpi.edu

Mark S. Shephard
Rensselaer Polytechnic Institute
Scientific Computation Research Center
shephard@rpi.edu

## MS68

### Shareable and Reproducible Cloud-based Experiments

Containerization has emerged as a systematic way of sharing experiments comprising of code, data, and environment. Containerization isolates dependencies of an experiment and allows the computational results to be regenerated. Several new advancements within containerization make it further easy to encapsulate applications and share lighter-weight containers. However, using containerization for cloud computing experiments requires further improvements both at the container runtime level and the infrastructure-level. In this presentation, we will first describe reproducible containers that encapsulate code, data and environment. We will then lay a vision for using containers as a dominant method for efficient sharing and improved reproducibility of cloud computing experiments. We advocate use of container-compliant cloud infrastructures, inclusion of performance profiles of the application or system architecture on which experiments were performed, and methods for statistical comparison across different container executions. We also outline challenges in the achieving this vision and propose existing solutions that can be adapted and propose new methods that can help with automation.

Tanu Malik
DePaul University
tanu.malik@depaul.edu

## MS68

### Towards Evidence-Based Best Practices for Reproducibility: A Software Engineering Research Perspective

In computational science and engineering (CSE) it is commonly understood that high-quality software and good development practices can help ensure reproducibility. At present, however, there is limited empirical support for how best to design for, implement, or maintain reproducibility over the course of the software lifecycle. As the scope, scale, and complexity of CSE software continues to grow, the state of practice must continue to evolve; looking to the horizon, emerging approaches to scientific computing (e.g. AI/ML-enabled applications and HPC-enabled edge computing) and software and hardware technologies (e.g. containers and novel accelerators) pose challenges and opportunities to improve how our community works with software. For that reason, there is an urgent need for (1) more research into practices, processes, and tools that would help scientific software teams seeking to build and

use reproducible software and (2) more efforts to socialize the successes of our community and promote the exchange of knowledge. In this talk, we will explore reproducibility from a software engineering perspective. That is, we will frame reproducibility as an engineering problem, where teams strategically invest in reproducibility alongside other quality goals and select practices that best serve their objectives. We will explore what we have learned from recent studies and offer actionable guidance for scientific software teams seeking to engineer and use reproducible software.

Reed Milewicz
Sandia National Laboratory
rmilewi@sandia.gov

Addi Malviya Thakur
Oak Ridge National Laboratory
malviyaa@ornl.gov

## MS68

### Ensuring Numerical Reproducibility in Scientific Codes Amidst the Challenges of Approximate Computing

Scientific computing is critical for the advancement of science and engineering, where numerical simulations are used for understanding complex phenomena, optimizing designs, and making predictions. With the slowdown of Moores law and the end of Dennard scaling, researchers are looking at other paradigms to improve the performance of scientific codes. One such promising approach is Approximate Computing, which involves trading-off accuracy for computational efficiency. While this approach offers significant benefits in terms of performance, it raises concerns about the reproducibility and correctness of scientific results. This talk explores the critical issue of numerical reproducibility in scientific codes. We will delve into the concepts of reproducibility, discussing what causes non-reproducibility and how to ensure the trustworthiness of computational results.We will discuss strategies for mitigating the impact of approximate computing on numerical reproducibility, including algorithmic modifications, reproducible linear algebra libraries, and error analysis.

Harshitha Menon
Lawrence Livermore National Laboratory
harshitha@llnl.gov

## MS68

### Transparency and FAIR are Steppingstones for Reproducibility in Science

The reproducibility crisis has not spared computational science, worsened by AI's growing role in research lifecycle. Scientific computing communities now employ intricate, data-heavy workflows that involve running multiple simulations concurrently, combining simulations with analysis, often employing machine learning, and orchestrating data-driven workflows with various task dependencies. These workflows pose numerous challenges to reproducibility. Merely packaging code into a container is insufficient; access to data and workflow scripts is necessary, as is documenting the execution graph where task order can impact reproducibility. Additionally, machine learning results have their own reproducibility constraints. Depending on the experiment's objective, reproducibility may vary along a spectrum. Characterizing the tradeoffs at the desired level of reproducibility and predicting the factors influencing variability of performance and results requires research in metrics, benchmarks, methods, and infrastructure. Transparency is a first step towards improved reproducibility in science, enabling re-use of data, the scrutiny of model parameters, and ultimately replication of experiments and results by teams other than the original developers. Embracing the FAIR (Findable, Accessible, Interoperable, Re-usable) principles and provenance is crucial for achieving computational reproducibility in data and workflows, fostering better scientific discovery reproducibility.

Line Pouchard
Brookhaven National Laboratory
pouchard@bnl.gov

## MS69

### Performance Engineering of the Navier-Stokes Finite Element Assembly of Alya on GPUs

This talk presents the measures taken to optimize an OpenACC GPU port of an existing, CPU friendly Fortran code base for a speedup of more than 50x on an NVIDIA A100 GPU. An analysis of the GPU and CPU performance of this initial OpenACC port identifies the reasons for sub optimal GPU performance. We believe that the anti patterns present in the original code are not uncommon in initial GPU ports using unified code bases, and that these measures are applicable to other code bases with similar development histories. The newly developed, GPU friendly code is beneficial for CPU performance as well, for a performance increase of more than 5x. The optimization subject is the assembly of the right-hand term in the incompressible flow module of the High-Performance Computational Mechanics code Alya, one of the two CFD codes in the Unified European Benchmark Suite.

Dominik Ernst
FAU university
dominik.ernst@fau.de

## MS69

### Efficient Algorithms for the Simulation of Incompressible Fluid Flow on Complex Geometries

Discontinuous Galerkin finite element methods are well suited to simulate transport-dominated problems on complex geometries, mainly due to their favorable dispersion and dissipation behavior and their geometric flexibility. In the context of incompressible fluid dynamics, we present DG methods designed with a focus on robustness and accuracy by a proper treatment of mass conservation and energy stability. An efficient implementation of such numerical methods on modern computer hardware is key to a successful realization in engineering applications, which are characterized by complicated geometries and under-resolved flow scenarios. We present the discontinuous Galerkin software project ExaDG, which realizes

high-order discretization methods for complex geometries through efficient matrix-free operator evaluation, fast iterative solvers and multigrid preconditioners, and efficient time integration methods. The utilized matrix-free algorithms are tuned towards optimal usage of memory bandwidth and caches on the node-level as well as parallel scalability on massively parallel systems. In this context, we discuss one of todays key challenges of CFD software on modern computer hardware (making progress in peak performance or memory bandwidth rather than in latency), namely the problem of minimizing the wall-time per time in the strong-scaling limit, which is vital in order to be able to simulate several convective time units for engineering applications.

Niklas Fehn, Martin Kronbichler
University of Augsburg
niklas.fehn@uni-a.de, martin.kronbichler@uni-a.de

## MS69

### Nonlinear FETI-DP Domain Decomposition Methods Combined with Deep Learning

In nonlinear-FETI-DP domain decomposition methods the choice of the nonlinear elimination set and of the coarse space have a huge impact on the nonlinear and linear convergence behavior. In this talk, we will show new results combining recently developed approaches for the adaptive choice of the nonlinear elimination set with adaptive coarse spaces. Additionally, we will discuss approaches to improve the computational efficiency and nonlinear convergence by enhancing Nonlinear-FETI-DP with techniques from machine learning.

Axel Klawonn
Universität zu Köln
axel.klawonn@uni-koeln.de

Martin Lanser
Universitaet zu Koeln
Mathematisches Institut
martin.lanser@uni-koeln.de

Janine Weber
Universität zu Köln
Department of Mathematics and Computer Science
janine.weber@uni-koeln.de

## MS69

### Efficient Implementations of High-Order Finite-Element Methods with Matrix-Free Operator Evaluation

This talk presents results from the collaborative research project PDExa, an initiative to advance algorithms for solving partial differential equations on exascale supercomputers. The main focus is on the development of efficient iterative solvers for challenging application problems in fluid dynamics discretized with high-order finite element methods. The application efficiency relies on three pillars. Fast matrix-free evaluation of the discretized PDE, computing the cell and face integrals on general curved meshes on the fly with sum-factorization techniques, implement the

action of linear or non-linear operators. In the talk, both classical $H^1$ and $L^2$ conforming methods will be considered and compared against $H(\text{div})$ conforming Raviart-Thomas operators for hyperbolic and elliptic terms in PDEs. For competitive iterative solvers, a second ingredient are preconditioners that balance low iteration counts with fast operator evaluation, for which either matrix-free and matrix-based ingredients can be attractive. Our project considers multigrid techniques, local solvers based on the fast diagonalization method and approximate incomplete matrix factorizations. As a third pillar, node-level performance engineering is employed, such as novel algorithms for higher data locality, aiming to use of the arithmetic capabilities of modern GPU and CPU hardware as efficiently as possible.

Martin Kronbichler
University of Augsburg
martin.kronbichler@uni-a.de

Ivo Dravins
Ruhr-University Bochum
ivo.dravins@rub.de

Niklas Fehn
University of Augsburg
niklas.fehn@uni-a.de

Marcel Koch
Karlsruhe Institute of Technology
marcel.koch@kit.edu

Katharina Kormann
Ruhr University Bochum
k.kormann@rub.de

## MS70

### Heterogeneity and HPC: The Messy Price for Performance

High Performance Computing (HPC) platforms continue to embrace diverse and novel hardware, all the while requiring that core software stacks remain functional, e.g., an implementation of the Message Passing Interface (MPI) or BLAS and LAPACK. This talk discusses some of the challenges that continued heterogeneity introduces both as a technical challenge for system administration and as an additional layer of complexity that users must become aware of.Unfortunately, this complexity spans all levels of the system, from administration to research and development and on to the user experience. In this talk, we will give some practical instances where challenges arise as well as discuss how co-design can be used to mitigate some of these risks and challenges.

James Elliot
Sandia National Laboratories
jjellio@sandia.gov

## MS70

### Scaling El Capitan: Application Preparations for

**DOE NNSAs First Exascale Computer**

El Capitan will be the DOE National Nuclear Security Administrations first exascale supercomputer. In collaboration with our vendor partners HPE and AMD through the El Capitan Center of Excellence (COE), we have spent several years preparing a variety of scientific applications such that they are scalable and performant on El Capitan on Day 1. In this talk, we will overview the application readiness efforts for El Capitan, discuss the modular software strategy employed by LLNL applications to ensure both performance and portability across modern architectures, provide early glimpses into the expected achievements that El Capitan will enable, and consider lessons learned and best practices developed by the El Capitan COE.

Judith Hill, Ramesh Pankajakshan
Lawrence Livermore National Laboratory
hill134@llnl.gov, pankajakshan1@llnl.gov

## MS70

**Porting Strategies for Targeting a GPU Supercomputer at CEA**

Modern supercomputers rely more and more on accelerators, such as GPUs, to provide the computing power required by simulation codes. CEA both hosts some of the most powerful supercomputers in the world, and develops numerous applications to simulate physics phenomena. With the advent of GPUs, CEA needs to adapt its codes to such architectures. Different strategies have been adopted depending on whether it is a new code or a legacy code, and on the complexity of the source code. We will present our effort to port CEAs codes to GPU machines. With heterogeneous parts having already be written in CUDA over the last years to target our own supercomputer, we evaluated the use of several more generic programming models such as CUDA, OpenMP target, Kokkos, SYCL, or through an adaptation of the ARCANE framework to GPUs. This presentation will focus on the lessons learned with targeting a generic programming model for GPUs.

Julien C. Jaeger
CEA
julien.jaeger@cea.fr

## MS70

**Towards Deeper Co-Design of Advanced Technologies**

Disruptive technology trends such as the slowing of Moore's law, open source and licensable hardware IP, fabless semiconductor firms, and Chiplets have motivated an exciting new approach to the co-design of Advanced Technology Systems at Los Alamos National Laboratory. This "Deeper" Co-Design will enable future technologies to be tailored to key aspects of our most challenging multiphysics problems in ways that were previously unattainable while building upon prior successes in Co-Design of systems such as Crossroads (Intel Xeon Max) and Venado (Nvidia Grace-Grace superchip and Grace-Hopper). Results of our work, including performance analysis on these three architectures and early results of hardware simulation

of alternative/future system designs will be presented.

Galen Shipman
Los Alamos National Laboratory
gshipman@lanl.gov

## MS71

**An Adaptive Algorithm Based on Partially Pivoted LU for Interpolative and CUR Decompositions**

Interpolative and CUR decompositions are special low-rank approximations in that they use original matrix rows/columns as bases. So they preserve properties of the original matrix (e.g., sparsity) and are easily interpretable. Many existing algorithms are based on column-pivoted QR factorization, which is challenging to parallelize due to the inherently sequential pivoting step. We present a new algorithm based on partially pivoted LU, which attains high performance on modern computing architectures such as multi-core CPUs and GPUs. We show that the new algorithm is fast, adaptive, and accurate.

Katherine J. Pearce
University of Texas at Austin
katherine.pearce@austin.utexas.edu

Chao Chen
North Carolina State University
chenchao.nk@gmail.com

Yijun Dong
New York University
yd1319@nyu.edu

Per-Gunnar Martinsson
University of Texas at Austin, U.S.
pgm@oden.utexas.edu

## MS71

**Robust Randomized Preconditioning for Kernel Ridge Regression**

We advocate two randomized preconditioning approaches for applying kernel ridge regression (KRR) to a moderate or large number of data points ($N \geq 10^4$). RPCholesky preconditioning is guaranteed to solve the exact KRR equations involving the $N \times N$ kernel matrix in just $\mathcal{O}(N^2)$ operations, assuming eigenvalue decay. KRILL preconditioning is guaranteed to solve the restricted KRR equations involving a $N \times k$ kernel submatrix in just $\mathcal{O}((N + k^2)k \log k)$ operations, with no assumptions on the kernel matrix or the regularization parameter. Experiments with dozens of data sets validate the effectiveness of RPCholesky and KRILL. Additionally, our theoretical analysis shows that RPCholesky and KRILL have stronger robustness properties compared to other commonly used preconditioners.

Mateo Diaz
Cornell University
mateodd@jhu.edu

Ethan N. Epperly
California Institute of Technology

eepperly@caltech.edu

Zachary Frangella
Stanford University, U.S.
zfran@stanford.edu

Joel A. Tropp
Caltech
jtropp@cms.caltech.edu

Robert Webber
California Institute of Technology
rwebber@caltech.edu

## MS71

### Robust Blockwise Random Pivoting: Fast and Accurate Adaptive Interpolative Decomposition

Interpolative decomposition (ID) aims to construct a low-rank approximation formed by a basis consisting of row/column skeletons in the original matrix and a corresponding interpolation matrix. This work explores fast and accurate ID algorithms from five essential perspectives for empirical performance: (a) skeleton complexity that measures the minimum possible ID rank for a given low-rank approximation error, (b) asymptotic complexity in FLOPs, (c) parallelizability of the computational bottleneck as matrix-matrix multiplications, (d) error-revealing property that enables automatic rank detection for given error tolerances without prior knowledge of target ranks, (e) ID-revealing property that ensures efficient construction of the optimal interpolation matrix after selecting the skeletons. While a broad spectrum of algorithms have been developed to optimize parts of the aforementioned perspectives, practical ID algorithms proficient in all perspectives remain absent. To fill in the gap, we introduce robust blockwise random pivoting (RBRP) that is parallelizable, error-revealing, and exact-ID-revealing, with comparable skeleton and asymptotic complexities to the best existing ID algorithms in practice. Through extensive numerical experiments on various synthetic and natural datasets, we empirically demonstrate the appealing performance of RBRP from the five perspectives above, as well as the robustness of RBRP to adversarial inputs.

Yijun Dong
New York University
Yijun Dong ¡yd1319@nyu.edu¿

Chao Chen
North Carolina State University
chenchao.nk@gmail.com

Per-Gunnar Martinsson
University of Texas at Austin, U.S.
pgm@oden.utexas.edu

Katherine J. Pearce
University of Texas at Austin

katherine.pearce@austin.utexas.edu

## MS71

### Novel Randomized Algorithms for QR with Column Pivoting, and Their Implementations in Randlapack

We present a pair of novel algorithms for QR decomposition with column pivoting. These algorithms use methods from RandNLA, in particular carefully using randomized sketching to accelerate both pivot decisions for the input matrix and the process of decomposing the pivoted matrix via Cholesky QR. The original algorithm, CQRRPT (pronounced "see-crypt"), is designed specifically for tall data matrices, while its successor, CQRRP, is applicable for matrices of any aspect ratio. We implement the algorithms in RandLAPACK by calling into RandBLAS and vendor-provided BLAS/LAPACK libraries. Experiments with these implementations were performed on an Intel Xeon Gold 6245 CPU and demonstrated order-of-magnitude speedups relative to LAPACKs standard function for QRCP, significant speedups over existing algorithms for QRCP and, in the case with CQRRPT, comparable performance to a specialized algorithm for unpivoted QR of tall matrices.

Max Melnichenko
University of Tennessee
mmelnic1@vols.utk.edu

Riley Murray
Sandia National Laboratories
rjmurr@sandia.gov

## MS72

### A Portable Multi-GPU Particle In Cell Electron Boltzmann Solver for Low-Temperature Plasmas

We study parallel particle-in-cell (PIC) methods for low-temperature plasmas (LTPs), which discretize kinetic formulations that capture the time evolution of the probability density function of particles as a function of position and velocity. We use a kinetic description for electrons and a fluid approximation for heavy species. In this paper, we focus on GPU acceleration of algorithms for velocity-space interactions and in particular, collisions of electrons with neutrals, ions, and electrons. Our work has two thrusts. The first is algorithmic exploration and analysis. The second is examining the viability of rapid-prototyping implementations using Python-based HPC tools, in particular PyKokkos. We discuss several common PIC kernels and present performance results on NVIDIA Volta V100 and AMD MI250X GPUs. Overall, the MI250X is slightly faster for most kernels but shows more sensitivity to register pressure. We also report scaling results for a distributed memory implementation on up to 16 MPI ranks.

James F. Almgren-Bell
UT Austin - Oden Institute
jalmgrenbell@utexas.edu

George Biros
University of Texas at Austin

Oden Institute
biros@oden.utexas.edu

## MS72

### Solving the Boltzmann Equation for Electron Kinetics Using Galerkin Approach

The collisional Boltzmann kinetic equations for low-temperature plasmas find important applications in semiconductor industry, materials science and many other applications in Science and Engineering. We present spatially homogeneous Boltzmann equation solver using Galerkin based deterministic and direct simulation Monte Carlo (DSMC) approaches. The developed Galerkin approach based deterministic solver supports user-specified multi-term expansion approximation for the distribution function with the support for electron-heavy and electron-electron Coulomb interactions. The above extends the traditional state-of-the-art two-term approximation (i.e., distribution function approximation using isotropic and anisotropic parts) used in the field. We provide detailed convergence studies, and cross-verification studies between the deterministic, DSMC and existing state-of-the-art Bolsig+ code.The presented multi-term Boltzmann solver is used to perform 1). Effects of the higher-order anisotropic correction terms on the traditional two-term approximation, 2). Relaxation time-scales for the perturbations in the higher-order anisotropic terms, and 3).Cross-verification of the higher-order correction terms using the DSMC approach.

Milinda Fernando
Oden Institute, UT Austin
milinda@oden.utexas.edu

George Biros
University of Texas at Austin
Oden Institute
biros@oden.utexas.edu

Robert D. Moser
University of Texas at Austin
rmoser@ices.utexas.edu

Laxminarayan Raja
Dept. of Aerospace Engineering and Engineering
Mechanics
The University of Texas at Austin
lraja@mail.utexas.edu

Todd A. Oliver
PECOS/ICES, The University of Texas at Austin
oliver@ices.utexas.edu

Philip Varghese
The University of Texas at Austin
varghese@mail.utexas.edu

## MS72

### Energy-Conserving and Sub-Cycled Electromagnetic Solvers for Scalable Particle-in-Cell Simulations of Low-Temperature Plasmas

The low temperature plasma group at the Princeton Plasma Physics Laboratory has been developing a 3-dimensional particle-in-cell kinetic code, LTP-PIC, for the study and rapid engineering prototyping of realistic industrial plasma devices. The code is designed from the ground up for scalability and performance, a critical requirement for the semiconductor industry and its needs for quick and iterative modeling tools. While the code has been optimized to take advantage of the impressive computing capabilities of modern GPU architectures, advanced algorithms are still necessary to achieve the fastest time to solution while maintaining high accuracy. In this talk, we will discuss the latest developments of the energy-conserving algorithm and sub-cycled electromagnetic solver implemented in LTP-PIC.

Stephane Ethier, Andrew Tasman Powis, Igor D.
Kaganovich
Princeton Plasma Physics Laboratory
ethier@pppl.gov, apowis@pppl.gov, ikaganov@pppl.gov

## MS73

### Scalable Bayesian Inference Using Integrated Nested Laplace Approximations

Integrated Nested Laplace Approximations (INLA) stand as a pivotal tool for performing approximate Bayesian inference, with broad applications in statistical modeling and various data-science fields. However, the current INLA implementation often hits computational and scalability ceilings, particularly when dealing with high-dimensional spatial-temporal models and intricate matrices. To overcome these limitations, this talk introduces a new versatile framework, designed specifically to leverage the scalability and processing power of manycore systems. This framework employs new numerical solver methods that take advantage of the dense block structure found in several data-science applications, thereby harmoniously blending the benefits of both dense and sparse computational techniques. By making optimal use of manycore architectures, the new framework effectively bypasses previous computational barriers, allowing for Bayesian inference on models with an extensive scope of latent parameters. We validate the efficacy of the new framework through comprehensive simulation tests and its application to real-world datasets. Our results confirm that the new framework significantly enhances computational performance while maintaining the rigor and accuracy of the inference, thereby pushing the envelope of what is feasible with the existing INLA implementation.

Esmail Abdul-Fattah
King Abdullah University of Science & Technology
(KAUST)
esmail.abdulfattah@kaust.edu.sa

## MS73

### oneMKL: Performance Portability from Libraries to Interfaces

Todays applications are more diverse than ever, as is the hardware used to run and accelerate them. oneAPI was

born to help overcome these challenges, by offering a unified, standards-based programming model that allows acceleration on various processing architectures. The oneAPI Math Kernel Library (oneMKL) element of the oneAPI specification defines SYCL-based APIs for a variety of mathematical routines, from dense and sparse linear algebra to discrete Fourier transforms to random number generators and more. The oneMKL interfaces project is an open-source implementation of the oneMKL specification and works with multiple hardware backends using device-specific libraries underneath to provide portable performance. Intel GPUs and CPUs benefit from the Intel oneMKL product, which supports both SYCL APIs as well as similar functionality for C and Fortran interfaces. This talk will present the recent advances in the Intel oneMKL product as well as the oneMKL interfaces project.

Sarah Knepper
Intel Corporation
sarah.knepper@intel.com

## MS73

### cuSOLVERMp: Challenges in Optimization Eigen-Value Problem for NVIDIA Based Superclusters

NVIDIA's cuSOLVERMp provides highly optimized routines for solving dense linear algebra problems on multinode systems. This presentation discusses the design principles and algorithmic choices that enable efficient support for heterogeneous systems, showcasing both scalability and performance results. Special attention is given to our state-of-the-art symmetric eigenvalue solver and demostrating how it unlocks the solution of massive problems in quantum chemistry and quantum mechanic research.

Samuel Rodriguez Bernabeu, Alexander Kalinkin
NVIDIA
srodriguezbe@nvidia.com, akalinkin@nvidia.com

## MS73

### Cholesky Factorization of H$^2$-matrices Without Trailing Matrix Dependencies

Structured low-rank matrices come in various formats such as BLR, BLR$^2$, HODLR, HSS, H-matrix, and H$^2$-matrix. Factorization of some of these matrices can be done in an inherently parallel fashion by separating the skeleton part from the redundant part of the dense blocks and factorizing only the redundant part at each level. For weak admissibility, this called the HSS-ULV factorization. However, for strong admissibility, the recompression of the fill-in blocks requires the shared basis to be updated, which in turn serializes the otherwise inherently parallel factorization. Our method aims to recover the inherent parallelism even for strong admissibility, by precomputing the fill-ins and including them in the shared basis before performing the actual factorization. This results in a highly parallel Cholesky factorization that can factorize structured low-rank matrices arising from 3-D geometry. The inherent parallelism allows us to use batch GPU kernels to achieve high performance. We also, develop a highly parallel method for the substitution part.

Rio Yokota
Tokyo Institute of Technology
rioyokota@gsic.titech.ac.jp

## MS74

### Scalable Low-Rank Optimization via Iteratively Reweighted Least Squares

Convex or non-convex matrix functions such as Schatten-p quasinorms or the positive log determinant have been successfully used as surrogates of the rank objective as an attractive reformulation to low-rank constraints, which in turn are ubiquitous in machine learning, computer vision, high-dimensional statistics and control. We review recent progress in the formulation and convergence analysis of Iteratively Reweighted Least Squares (IRLS), which has emerged as a uniquely suitable algorithmic framework to scalably optimize such rank surrogates. Furthermore, we present a highly parallelizable implementation of our approach that is tailored to low-rank matrix completion problem, which is widely used in recommender systems.

Christian Kümmerle, Nicholas Cassarino, Tyler Allen
University of North Carolina at Charlotte
kuemmerle@uncc.edu, ncassari@charlotte.edu, t.allen@charlotte.edu

## MS74

### Decentralized Algorithms for Spatially Distributed Systems

In a centralized system, all processing and decision-making is handled by a single entity, which can become a bottleneck as the system grows in size and complexity. Decentralized algorithms distribute the processing load across multiple nodes, allowing the system to scale much more effectively. Many decentralized algorithms distribute the global objective function (usually the sum of many local objective functions) across multiple nodes such that each node only handle its own local objective function. The state variable is copied to each node and communication between neighboring nodes helps them reach a consensus in the convergence of iterations. However, the computational cost for each agent can still be high if the common state variable has a large dimension. To address this, we would further divide the global state variable into multiple local state variables, so that each node only handles a few components of the global state variable. In particular, we will analyze its performance on spatially distributed systems.

Guohui Song
Department of Mathematics and Statistics
Old Dominion University
gsong@odu.edu

## MS74

### Distributed Learning for Kernel Mode-Based Regression

We propose a parametric kernel mode-based regression built on mode value, which can achieve robust and efficient estimators when the data have outliers or heavy-tailed distributions. We show that the resultant estimators can arrive at the highest asymptotic breakdown point of 0.5.

We then utilize such a regression for massive datasets by combining it with the distributed statistical learning technique, which can greatly reduce the required amount of primary memory while simultaneously incorporating heterogeneity into the estimation procedure. By approximating the local kernel objective function using a least squares format, we are able to preserve compact statistics for each worker machine and employ them to rebuild the estimate of the entire dataset with asymptotically minimal approximation error. With the help of a Gaussian kernel, an iteration algorithm built on expectation-maximization procedure is introduced, which could substantially lessen the computational burden. The asymptotic properties of the developed mode-based estimators are established, where we prove that the suggested estimator for massive datasets is statistically as efficient as the global mode-based estimator using the full dataset. We further conduct a shrinkage estimation based on the local quadratic approximation and demonstrate that the resulting estimator has the oracle property. The finite sample performance of the developed method is illustrated using simulations as well as real data analysis.

Tao Wang
Department of Economics
University of Victoria
taow@uvic.ca

### MS74

### Classification of Unbounded Data by Gaussian Mixture Models via Deep ReLU Networks

This paper studies the binary classification of unbounded data from $\mathbb{R}^d$ generated under Gaussian Mixture Models (GMMs) using deep ReLU neural networks. We obtain for the first time non-asymptotic upper bounds and convergence rates of the excess risk (excess misclassification error) for the classification without restrictions on model parameters. The convergence rates we derive do not depend on dimension $d$, demonstrating that deep ReLU networks can overcome the curse of dimensionality in classification. While the majority of existing generalization analysis of classification algorithms relies on a bounded domain, we consider an unbounded domain by leveraging the analyticity and fast decay of Gaussian distributions. To facilitate our analysis, we give a novel approximation error bound for general analytic functions using ReLU networks, which may be of independent interest. Gaussian distributions can be adopted nicely to model data arising in applications, e.g., speeches, images, and texts; our results provide a theoretical verification of the observed efficiency of deep neural networks in practical classification problems.

Tian-Yi Zhou, Xiaoming Huo
Georgia Institute of Technology
tzhou306@gatech.edu, huo@isye.gatech.edu

### MS75

### Manipulating Sparse Computations Using Dense Polyhedra: Challenges and Opportunities

Sparse computations typically operate on input data and their coordinates on an integer grid. These non-zero coordinates can be compressed using specific sparse formats, such as CSR or CSC. We have shown these non-zero coordinates can also be described using union of integer polyhedra, combined with integer lattices, to deliver optimized computation kernels that are specialized to a particular sparsity structure and free of any indirect array access, suitable for aggressive vectorization. In this talk we will present our approaches to manipulating sparse structures using union of polyhedra and generating efficient pattern-specific SIMD code; and recent results in developing a new sparse format based on polyhedra supporting a rich set of sparse computations kernels.

Louis-Noël Pouchet
Colorado State University
pouchet@colostate.edu

### MS75

### Towards Reverse Mode Automatic Differentiation of Kokkos-Based Code Using the Llvm Compiler Infrastructure

The Kokkos programming model is a production level solution for writing performance portable modern C++ applications in a hardware agnostic way. It has been used in multiple applications and simulation software to model molecular dynamics, solid mechanics, fluid mechanics, combustion phenomena, and plasma simulations among others. In most of those applications, the simulation code needs to compute some derivatives of some quantity of interest or of some residual with respect to parameters or solution state. The templated arguments of Kokkos view, the multidimensional arrays of the Kokkos programming model, have been used to implement forward Automatic Differentiation (AD) using operator overloading. While this approach is performance portable and very efficient when the code compute the derivative of a lot of quantities with respect to a few variables at once, it is not as efficient to compute the derivatives of a few quantities with respect to a lot of variables. That second case is however extremely important to solve optimization problems with a larger number of parameters or to train a neural network with a lot of neurons for instance. Instead of using a forward AD approach, these situations would benefit of a reverse AD approach. In this talk, we will discuss the usage of the LLVM compiler infrastructure to compute the derivatives of a few quantities of interest with respect to a lot of parameters using a source code transformation reverse mode AD.

Sivasankaran Rajamanickam, Kim Liegeois, Brian Kelley
Sandia National Laboratories
srajama@sandia.gov, knliege@sandia.gov,
bmkelle@sandia.gov

Eric Phipps
Sandia National Laboratories
Center for Computing Research

etphipp@sandia.gov

## MS75

### Tensql: SQL Database Using Graphblas

Relational Database Management Systems (RDBMS) have been the most prominent form of database in the world for several decades. While relational databases are often applied within high-frequency/low-volume transactional applications such as website backends, the poor performance of relational databases on low-frequency/high-volume queries often precludes their application to big data analysis fields like graph analytics. This work explores the construction of an RDBMS solution that uses the GraphBLAS API to execute Structured Query Language (SQL) in an effort to improve performance on high-volume queries. Tables are redefined to be collections of sparse scalars, vectors, matrices, and more generally sparse tensors. The explicit values (nonzeros) in these sparse tensors define the rows and NULL values within the tables. A prototype database called TenSQL was constructed and evaluated against several SQL implementations including PostgreSQL. Preliminary results comparing the performance on queries common in graph analysis applications offer performance improvements as high as 1,400x over PostgreSQL for moderately sized datasets when returning results in a columnar format.

Miheer Vaidya, Ponnuswamy Sadayappan
University of Utah
m.vaidya@utah.edu, saday@cs.utah.edu

Jon P. Roose, Sivasankaran Rajamanickam
Sandia National Laboratories
jproose@sandia.gov, srajama@sandia.gov

## MS75

### Distributed Sparse Computing in Python

The sparse module of the popular SciPy Python library is widely used across applications in scientific computing, data analysis, and machine learning. The standard implementation of SciPy is restricted to a single CPU and cannot take advantage of modern distributed and accelerated computing resources. We introduce Legate Sparse, a system that transparently distributes and accelerates unmodified sparse matrix-based SciPy programs across clusters of CPUs and GPUs, and composes with cuNumeric, a distributed NumPy library. Programs written in Legate Sparse and cuNumeric perform competitively with low-level systems like PETSc. In this talk, Ill discuss one aspect of the design and implementation of Legate Sparse called DISTAL, a compiler for sparse tensor algebra that targets distributed machines, that we used to generate a large portion of Legate Sparse. DISTAL separates descriptions of tensor algebra expressions, sparse data structures, data distribution and computation distribution, enabling the distributed execution of sparse tensor algebra expressions with a variety of sparse data structures and data distributions. DISTAL generates code that is competitive with hand-written implementations, and outperforms general interpretation based systems by one to two orders of magnitude

Rohan Yadav
Stanford University
rohany@stanford.edu

Wonchan Lee, Melih Elibol, Manolis Papadakis, Taylor Lee Patti, Michael Garland, Michael Bauer
NVIDIA
wonchanl@nvidia.com, melibol@nvidia.com, mpapadakis@nvidia.com, tpatti@nvidia.com, mgarland@nvidia.com, mbauer@nvidia.com

Alex Aiken, Fredrik Kjolstad
Stanford University
aaiken@stanford.edu, kjolstad@cs.stanford.edu

## MS76

### Learning Closure Relations Using Differentiable Programming

The continuous flow or 'transport' of a macroscopic system of particles is a high dimensional problem and therefore often solved using reduced order models. This necessarily introduces unknown closure relations into these models. In this work, we present a machine learning approach to finding accurate closure relations utilising differentiable programming. As a case study, we consider the transport of photons and use a literature radiation transport test problem as a training dataset. We present novel closures for a number of reduced order models. We evaluate the improvement of the machine-learnt closures over those from the literature.

Aidan Crilly
I-X Centre for AI In Science, Imperial College London
a.crilly16@imperial.ac.uk

## MS76

### Integrating Checkpoint Schedules with an Algorithmic Differentiation

In this work, we aim to present the conjunction of a checkpoint schedule package and the differentiation algorithm pyadjoint, a Python library capable of generating automated adjoints compatible with PDE frameworks Firedrake or FEniCS. The checkpoint schedule package offers checkpointing algorithmics employed to avoid the storage of the entire forward state in preparation for the adjoint computation, which is particularly critical for large-scale problems requiring significant memory usage. Accordingly, the checkpoint schedules and the pyadjoint integration is essential since this algorithmic differentiation generates automated adjoints for diverse time-dependent physical models tackled by Firedrake and FEniCS users. We focus on presenting the design of the checkpoint schedules and pyadjoint integration, as well as the numerical simulations obtained using the checkpointing algorithms generated by the checkpoint schedule package.

Daiane I. Dolci
University of São Paulo
d.dolci@imperial.ac.uk

David Ham
Department of Mathematics
Imperial College London
david.ham@imperial.ac.uk

## MS76

### Modeling for Inversion in Geophysics with Devito

Seismic inversion, and more generally geophysical exploration, aims at better understanding the Earth's subsurface. A plethora of challenges across a range of scientific disciplines ensure this is a challenging problem. In this work we will discuss how the Devito package has been developed and optimised to address some of the challenges present within geophysical inversion workflows. Devito is a Python based domain specific language and compiler for expressing the solutions of PDE boundary value problems in terms of finite differences. From inception, one overarching goal has been that the generated code is optimal within the context of seismic inversion problems: that is, the auto-generated computational kernels match or beat hand-optimized codes for forward wavefield and adjoint computation. This work will provide an overview of Devito's application in this area and also discuss the roadmap to broaden the applicability of Devito to additional areas such as CFD.

Rhodri Nelson, George Bisbas
Imperial College London
rnelson@ic.ac.uk, g.bisbas18@imperial.ac.uk

Matthias Louboutin
Georgia Institute of Technology
mathias.louboutin@gmail.com

Fabio Luporini
Devito Codes
fabio@devitocodes.com

Gerard Gorman
Imperial College London
g.gorman@imperial.ac.uk

## MS76

### The Lynchpin of the Adjoint

In the world of adjointing, people tend to fall into two tribes. The Team Numerical says the only true path is to adjoint discretised equations. They say "It's the way to be sure you only compute the gradient of what you start with." On the other hand, the Mathematical Squad says to take the analytical adjoint of the original system of equations and discretise the result later, saying, "It's the abstractly elegant way'. Empirically, it seems (based on much hearsay) that Team Numerical is winning along many fronts. However, there isn't much to reconcile the two camps based on first principles. This talk will present a unified approach to the two views. The key is the residual space of the approximated differential equation. Said another way, there is an exact analytical form that expresses the discrete adjoint. And considering the original differential system in an algorithm-agnostic fashion naturally points out the correct basis to discretise an adjoint

system. The residual space is the lynchpin of the adjoint.

Geoff Vasil
School of Mathematics, University of Edinburgh
gvasil@ed.au.uk

## PP1

### IPPL: A Massively Parallel, Performance Portable C++ Library for Particle-Mesh Methods - Structure, Use Cases and Implementation Details

We present implementation details and use cases of the Independent Parallel Particle Layer (IPPL), a performance portable C++ library for particle-in-cell methods. IPPL makes use of Kokkos (a performance portability abstraction layer), HeFFTe (a library for large scale FFTs), expression templates and MPI (Message Passing Interface) to deliver a portable, massively parallel toolkit for particle-mesh methods. IPPL supports simulations in one to six dimensions, mixed precision, and asynchronous execution in different execution spaces (e.g. CPUs and GPUs). We showcase the easy to use library, the performance of the latest version of IPPL using examples from charged particle dynamics (https://arxiv.org/abs/2205.11052) on state-of-the-art high-performance computing resources, such as Perlmutter and Frontier.

Andreas Adelmann
Paul Scherrer Institut
andreas.adelmann@psi.ch

Antoine Cerfon
Type One Energy
antoine.cerfon@typeoneenergy.com

Matthias Frey
University of St Andrews
mf248@st-andrews.ac.uk

Sonali Mayani
Paul Scherrer Institut
sonali.mayani@psi.ch

Veronica Montanaro
ETH
vmontanaro@student.ethz.ch

Sriramkrishnan Muralikrishnan
Juelich Supercomputing Centre
Forschungszentrum Juelich GmbH
s.muralikrishnan@fz-juelich.de

Alessandro Vinciguerra
ETH Zurich
vincigua@student.ethz.ch

## PP1

### One-Directional Mesh Overset Using a Parallel Distributed Forest of Octrees

Many physical systems cannot be modeled properly by a single mesh, e.g. because the domain contains holes, or several physics layers interact. A known natural solution to

this is the use of several overlapping meshes covering different parts of the domain. To obtain meaningful results from this mesh overset approach, the meshes have to yield a consistent solution wherever they overlap. To achieve this they need to exchange interpolated solution data. We present our current state of research, an algorithm to handle the overset between two parallel-distributed meshes. One mesh queries and receives the interpolated data. Its cells are discretized by query points (e.g. stemming from quadrature), thereby it may be of almost arbitrary type and structure. The second mesh is a forest of octrees, which provides the solution data. The forest-of-octrees structure allows for efficient searching of the query points in the local part as well as in the global partition of the mesh. By applying smooth, invertible mappings to the individual trees of the forest, one can model complex geometries while preserving the possibility to search points in the axis-aligned tree structure taken for reference. The query points are sent between the two meshes by non-blocking point-to-point MPI communication, which allows to overlap communication with computation. We outline the sub-steps of our new generic mesh overset algorithm and present numerical results showcasing its flexibility and scalability.

Hannes Brandt
University of Bonn
brandt@ins.uni-bonn.de

Carsten Burstedde
Universität Bonn
burstedde@ins.uni-bonn.de

**PP1**

**Neural Network Solver for Transient State Incompressible Flow on Arbitrarily-Shaped Geometries Defined Using Bezier Curves**

Models for simulating turbulent fluid flow have many industrial applications. Current methods for modeling turbulent fluids are accurate but slow. As a result, much research has been published on neural networks that can predict fluid flow with inference times that are orders of magnitude faster than traditional methods. To enable the neural network to predict the fluid flow around any arbitrarily-shaped obstruction, the prevailing technique is to input the simulation grid into the neural network in the form of a point cloud or image. In this research, we propose a simpler approach. Instead of a point cloud, we define the boundary of an arbitrary obstruction using a handful of points. All intermediate points will then be interpolated using Bezier curves. Using fewer points allows for improved training time, inference time, and neural network size. Further, by applying parallel processing techniques, we optimized both speed and accuracy. During the data generation step, parallelizing independent loops reduced our time cost by 21.0%; ensemble stacking reduced inference time by 53.1%; finally, taking the average result of the ensemble decreased the error by 45.8% (versus any one neural network). After 90 epochs, our ensemble was able to predict transient-state turbulent fluid flow with a mean squared error of 1.162E-05. Our proposed technique is applicable to modeling fluid flow over 3D topographic maps, which are suited for interpolation using curves.

Jadey Chen
Gabrielino High School
jadeyc124@gmail.com

**PP1**

**An Analysis on Multigrid Reduction in Time Convergence of Time Averaged Quantities in High Speed Flow**

The Multigrid Reduction in Time (MGRIT) Algorithm [Friedhoff et. al., A Multigrid in Time Algorithm for Solving Evolution Equations in Parallel, 2013], a time domain analogue to spatial Multigrid Reduction, has been shown to offer convergence speedup in the solution of uncomplicated parabolic cases [Falgout et. al., Parallel Time Integration with Multigrid, 2013]. We present an analysis of the extent to which MGRIT can be extended to a high speed, time dependent, fluid flow simulation. Specifically, we test the hypothesis that certain time averaged statistics on the flow can be computed faster with MGRIT. To do so, we measure fluctuating quantities relevant to a two dimensional flow based on the Euler equations around a cylindrical object using the software package XBraid [XBraid: Parallel multigrid in time. http://llnl.gov/casc/xbraid] to handle the MGRIT algorithm, and deal.II [Arndt et. al., The deal.II Library, Version 9.5, 2023] for the spatial discretization.

Jerett C. Cherry
Colorado State University
jerett97@colostate.edu

**PP1**

**Scalable Hierarchical Approximation of Dense Kernel Matrices**

Data-driven approaches in areas like machine learning and deep learning have led to rapid expansion in data sizes. Dense matrices of size $N \times N$ have extremely high memory requirements of $O(N^2)$ and vital methods like eigendecompositions become computationally infeasible. We apply a novel method (GOFMM) to compute a hierarchical approximation of dense matrices found in manifold learning algorithms such as diffusion maps. GOFMM "compresses" the matrix such that memory requirements and computational cost of operations such as the matrix-vector multiplication reduce from $O(N^2)$ to $O(NlogN)$. We then test the accuracy of our method by calculating error norms of eigenpairs obtained before and after the hierarchical approximation. We also test for scalability of our approach while increasing matrix sizes upto $16K \times 16K$ on multiple nodes using MPI parallelism. Our work-in-progress extends the application of GOFMM for dense kernel matrices obtained in Convolutional Neural networks.

Keerthi Gaddameedi
Technical university of munich
keerthi.gaddameedi@tum.de

Severin Reiz, Felix Dietrich, Tobias Neckel
Technical University of Munich
s.reiz@tum.de, felix.dietrich@tum.de, neckel@in.tum.de

Hans Joachim Bungartz
Technical University Munich
bungartz@tum.de

## PP1

### A Minimal, Serial-Equivalent Format for Parallel I/O and Its Applications

Large-scale simulations execute as parallel jobs that partition data among multiple processes. If each process writes its local data into a separate file, we encounter practical limitations because the writing partition is required for reading the data files. Therefore, we specify an application- and code-agnostic file-oriented data format suitable for parallel, partition-independent disk I/O. Here, a partition refers to a disjoint and ordered distribution of the data elements between one or more processes. The format is designed such that the file contents are invariant under linear (i.e. unpermuted) parallel repartitioning of the data prior to writing. The file contents are indistinguishable from writing in serial. In the same vein, the file can be read on any number of processes that agree on a partition for the number of elements stored. Out of necessity the sheer size of the data needs to be taken into account as well. Therefore, we add an optional convention to implement transparent, lossless, per-element data compression. The compressed data and metadata is layered inside ordinary format elements and likewise partition-independent. We refer to this format as scda. The main purpose of the format is to abstract any parallelism and provide sufficient structure as a foundation for generic and flexible archival and checkpoint/restart. We provide parallel simulation checkpoint/restart examples built on top of the adaptive forest-of-octrees software library p4est.

Tim Griesbach
Universität Bonn, Germany
Institute for Numerical Simulation
tim.griesbach@uni-bonn.de

Carsten Burstedde
Universität Bonn
burstedde@ins.uni-bonn.de

## PP1

### P4est3: Performance and Interface Extensions for the Parallel AMR Library P4est

Solving modern engineering and scientific problems requires high accuracy calculations at minimized computational complexity. We consider adaptive mesh refinement (AMR) as a key technique for simulations requiring time-dependent and multiscale features. p4est is parallel-distributed, scalable software library for managing AMR. It empowers efficient partitioning, neighbor finding, remote object search, mesh entity iteration, and more. This poster showcases advancements in p4est through an experimental extension pack. First, it introduces a general interface for handling diverse cells representations. This extension adds two more implementations for encoding mesh primitives: as a linear index and as an entity within extended AVX/SSE CPU registers. The next enhancement exploits MPI-3 shared memory. This eliminates redundancy of quadrant and metadata storage, allows for a message-free partitioning algorithm and eliminates the need for a ghost layer structure for mesh iteration within a single shared memory node. The third improvement is parallel algorithms optimization, featuring faster top-down forest creation, generic iteration interface without 2:1 mesh-balance requirement and more general refining and coarsening. Additionally, it introduces RAII and COW approaches to a forest lifecycle. In conclusion, we demonstrate how simulation pipelines differ for p4est with and without the proposed extensions, and we compare the performance in various configurations.

Mikhail Kirilin, Carsten Burstedde
Universität Bonn
kirilin@ins.uni-bonn.de, burstedde@ins.uni-bonn.de

## PP1

### Enhancing Weakly Compressible Smoothed Particle Hydrodynamics with Lagrangian Parallelism for Free Surface Flow Simulation

In this study, we present a novel approach to accelerating weakly compressible smoothed particle hydrodynamics (WCSPH) simulations on CPU platforms. Our research focuses on the development of a WCSPH code optimized through the innovative integration of Lagrangian parallelism. This code is specifically designed for the numerical simulation of free surface flow, a critical domain in fluid dynamics. The Lagrangian parallelism technique employed in our work revolutionizes the distribution of computational tasks across processors. It is tailored to align seamlessly with the Lagrangian nature inherent in fluid simulations. To validate the effectiveness of our method, we conducted extensive testing on two key test cases: the stretching drop of water and the dam-break scenario. These cases were respectively compared against an analytical solution and an established experimental benchmark. Our results demonstrate a remarkable equilibrium between computational precision and performance. By leveraging Lagrangian parallelism, we have not only enhanced the accuracy of WCSPH simulations but also achieved notable improvements in computational efficiency. This research contributes significantly to the field of parallel processing for scientific computing, offering a promising avenue for advancing free surface flow simulations.

Mohamed Labadi
University of Chlef
Frontiers Labs
mohamed.labadi2@gmail.com

Abdelkader Krimi
Ecole Polytechnique Montreal, Montreal-Canada
abdelkader.krimi@polymtl.ca

Samir Abdelmalek
University of Medea
Complex Systems Labs
abdelmalek.samir@univ-medea.dz

## PP1

### Towards Automatic Adjoint Differentiation via

## Symbolic Computation

Adjoint-based sensitivity analysis is an invaluable tool for forecasting and optimisation frameworks, though its implementation is limited in large scale applications such as wind farm optimisation. Symbolic computation, where optimised code is automatically generated from a high-level problem definition, offers an efficient method for building adjoint models for an extended range of computations. An example is the adjoint subpackage integrated within the automated finite element PDE solver Firedrake. To tackle large scale sensitivity problems with finite difference computations, we aim to extend the domain specific language Devito with similar adjoint capabilities. Devito is an industry standard Python package used to generate optimised stencil computations from high-level symbolic abstractions over multiple computer architectures. First, the implementation of implicit solver functionality into Devito is necessary for the applications of interest, such as adjoint-based wind farm optimisation. This poster will display an early prototype interfacing with the PETSc infrastructure, wrapping matrix-free routines around the operator application in Devito in order to automatically generate and run the underlying linear solvers. Whilst early in development, this approach will enable Devito to tackle a wider range of challenges in the field of high-performance computing.

Zoe Leibowitz
Imperial College London
zoe.leibowitz21@imperial.ac.uk

Rhodri Nelson
Imperial College, United Kingdom
rhodri.nelson@imperial.ac.uk

Fabio Luporini
Devito Codes
fabio@devitocodes.com

Mathias Louboutin
Georgia Institute of Technology
mathias.louboutin@gmail.com

Gerard Gorman
Imperial College London
g.gorman@imperial.ac.uk

Matthew G. Knepley
University at Buffalo
Department of Computer Science and Engineering
knepley@gmail.com

## PP1

## Consistent Coupling of Multiscale MPAS and ROMS Ocean Models

Numerical resolution plays a crucial role in maintaining model accuracy but is directly tied to significant computational expenses, which pose limitations on the development of high-fidelity global Earth system models (ESMs) for long-term climate simulations. Despite these challenges, there is a strong interest in understanding the effects of coastal and open-ocean dynamics on climate, necessitating considerably higher resolutions. We introduce an infrastructure designed to advance our understanding of ocean processes, employing a high-resolution Regional Ocean Modeling System (ROMS) alongside a global Model for Prediction Across Scales-Ocean (MPAS-O) for regions of interest, in order to enhance the fidelity of nonlinear ocean dynamics. This multiscale ocean model necessitates seamless one-way and two-way coupling and field transfers to efficiently assess climate variability while upholding model precision. We illustrate the improved representation of small-scale processes using a flexible MOAB library-based infrastructure.

Vijay S. Mahadevan, Iulian Grindeanu
Argonne National Laboratory
mahadevan@anl.gov, iulian@anl.gov

Robert Hetland
Pacific Northwest National Laboratory
robert.hetland@pnnl.gov

## PP1

## Broadening Participation in Doe High-Performance Computing

In high-performance computing (HPC), traditionally underrepresented groups engage at significantly lower rates than their representation in the general population. As a unique multilab partnership across DOE computing sciences, the Exascale Computing Project (ECP) has been uniquely positioned to contribute to closing this gap. The ECP Broadening Participation Initiative has established a sustainable plan to recruit and retain a diverse HPC workforce by fostering a supportive and inclusive culture within the computing sciences at DOE national laboratories. Our approach has three complementary thrusts that leverage existing efforts on workforce development in DOE national laboratories, computing facilities, and the HPC computational science community: 1) the HPC Workforce Development and Retention Action Group, a community-based working group of DOE lab staff and collaborators, who exchange ideas and best practices on building and sustaining an inclusive HPC workforce; 2) the Sustainable Research Pathways Program, a multi-lab internship and workforce development program that pairs national lab staff with students from underrepresented groups (and faculty working with them) to partner on research projects; and 3) the Introduction to HPC Bootcamp Program, which provides introductory HPC training to undergraduates and early-career graduate students, both in computing and domain sciences through energy justice project-based learning activities.

Daniel Martin
Lawrence Berkeley National Laboratory
dfmartin@lbl.gov

Paige Kinsely
Argonne National Laboratory
pkinsley@anl.gov

Mary Ann E. Leung
Sustainable Horizons Institute
mleung@shinstitute.org

Lois Curfman McInnes
Argonne National Laboratory
curfman@anl.gov

Suzanne Parete-Koon
Oak Ridge National Laboratory, U.S.
paretekoonst@ornl.gov

Sreeranjani Ramprakash
Argonne National Laboratory, U.S.
ramprakash@anl.gov

## PP1

### SIAG/Supercomputing Initiatives: Raising Awareness of HPC Opportunities and Impact Growing the HPC Community

The SIAM Activity Group on Supercomputing (SIAG/SC, https://siag-sc.org) provides a forum for computational mathematicians, computer scientists, computer architects, and computational scientists to exchange ideas on mathematical algorithms and computer architecture needed for high-performance computer systems. SIAG/SC promotes the exchange of ideas by focusing on the interplay of analytical methods, numerical analysis, and efficient computation. This poster provides an overview of SIAG/SC initiatives that aim to raise awareness of opportunities and impact in high-performance computing (HPC) and grow the HPC community. A focus is the Supercomputing Spotlights webinar series, featuring short presentations that highlight the impact and successes of HPC throughout our world. Presentations, emphasizing achievements and opportunities in HPC, are intended for the broad international community, especially students and newcomers to the field. We welcome your ideas and contributions Join us! SIAG/SC Officers 2022-2023: Lois Curfman McInnes (chair), Hatem Ltaief (vice chair), Michael Bader (program director), Rio Yokota (secretary)

Lois C. Mcinnes
Mathematics and Computer Science Division
Argonne National Laboratory
curfman@anl.gov

Hatem Ltaief
KAUST, Saudia Arabia
hatem.ltaief@kaust.edu.sa

Michael Bader
Technical University of Munich
bader@in.tum.de

Rio Yokota
Tokyo Institute of Technology
rioyokota@rio.gsic.titech.ac.jp

## PP1

### AutoPas: Dynamic Algorithm Selection in Molecular Dynamics for Optimal Time and Energy

The large computational cost of force calculations within Molecular Dynamics has led to the development of specialist algorithms such as Linked Cells or Verlet Lists, each with multiple parallelisation schemes. There is, however, no single best algorithm for all scenarios and architectures, and the best algorithm can vary across the domain of a simulation and can change over time. In this work, we will present AutoPas: a black-box particle simulation library that aims to select the fastest or most energy-efficient algorithm for a given simulation, including using different algorithmic choices for different regions of the domain and dynamically updating the algorithmic choices as the simulation changes [F. Gratl et al., 2022, N ways to simulate short-range particle systems: Automated algorithm selection with the node-level library AutoPas]. In particular, we will discuss the impact of simulation-specific properties, such as density and homogeneity, on the performance of the algorithms and how machine learning and expert knowledge can make use of such properties to ensure optimal algorithm selections with minimal overhead.

Samuel J. Newcome, Fabio A. Gratl, Markus
Muehlhaeusser
Technical University of Munich
samuel.newcome@tum.de, f.gratl@tum.de,
markus.muehlhaeusser@tum.de

Philipp Neumann
Helmut Schmidt University, Germany
philipp.neumann@hsu-hh.de

Hans-Joachim Bungartz
Technical University of Munich
bungartz@cit.tum.de

## PP1

### An Analytical Diagonalization Technique for Approximating the Spectral Fractional Laplacian

Problems involving non-local operators have recently attracted increasing interest in many diverse fields. However, non-locality necessarily increases the computational complexity to approximate solutions to these problems. We study the spectral fractional Laplacian $(-\Delta)^s u = f$ in a bounded domain $\Omega \subset \mathbb{R}^d$. Previous works have used the Caffarelli-Silvestre extension to convert the fractional Laplacian into a Dirichlet-to-Neumann mapping in $\mathbb{R}^{d+1}_+$. A diagonalization scheme is used to reduce the computational complexity by exposing the inherent parallelizability of the method. We refine the diagonalization scheme by proposing an analytic approach to compute the eigenpairs of the eigenvalue problem in the extended dimension, avoiding the numerical instability in approximating the eigenpairs with a finite element method. We demonstrate that this new analytical approach is related to certain quadrature schemes used to approximate the spectral fractional Laplacian. We further show that this novel algorithm maintains exponential convergence. Numerical examples in two dimension demonstrate the performance of the method.

Shane E. Sawyer
University of Tennessee, Knoxville
sjw355@vols.utk.edu

Abner J. Salgado
Department of Mathematics

University of Tennessee
asalgad1@utk.edu

## PP1

**Using a Hybrid Gpu/cpu Parallel Bayesian Inversion to Determine the Physics and Rupture Dynamics Governing the 2004 Parkfield, Ca, Earthquake**

Earthquakes are caused by frictional failure on weak faults within the Earth's crust. Seismologists routinely determine the kinematics of earthquake sources by inverting surface observations but the results are inherently non-unique. 40+ seismic and 13 GPS stations recorded similarly large co- & postseismic slip of 2004 Mw 6.0 Parkfield earthquake. We employ a hybrid GPU/CPU parallel approach to invert for the underlying stress and friction parameters. This is a highly non-linear, very high-dimensional problem (without the possibility of using gradients), which involves a computationally expensive dynamic forward problem consisting of two phases and a total of 1100 model parameters. The coseismic phase is described by the hyperbolic wave equation and the postseismic phase by the elliptic problem of linear elasticity, both coupled to an internal rate-and-state friction boundary condition. We constrain the inversion with seismic and GPS data. Our hybrid approach, in which we compute the coseismic phase on the GPU and the postseismic phase on the CPU, allows us to use 3 MPI processes per GPU. Our cluster with 8 Nvidia RTX A5000 GPUs constantly ran for more than 8 months to produce several millions of forward solutions. Our final model ensemble illuminates the dynamics of the 2004 Parkfield earthquake with unprecedented accuracy and provides estimates of stress drop, fracture energy, and radiation efficiency, which are usually hard to infer from pure surface observations.

Nico Schliwa
Ludwig-Maximilians-Universität München
nico.schliwa@geophysik.uni-muenchen.de

Alice-Agnes Gabriel
UC San Diego
algabriel@ucsd.edu

Jan Premus
Charles University Prague
janpremus@seznam.cz

František Gallovic
Charles University, Prague
frantisek.gallovic@matfyz.cuni.cz

## PP1

**A Fully Adaptive Iterative Scheme for Solving Bilevel Variational Inequality Problems.**

The study of Bilevel Variational Inequality Problems (BVIP) has long been a subject of intense research interest, owing to the diverse applications in various fields such as science, engineering, medicine, cryptography, image and signal processing, and optimal control. In this presentation, we give overview of some methods proposed over time to solve variational inequality problems including bilevel case. In particular, this talk will introduce our recently proposed iterative scheme for solving BVIP involving monotone operators. The presence of the inertial parameter and more efficient step-size make the proposed algorithm a very robust scheme. Furthermore, our iterative scheme involves a single projection onto half-space which contributes to the reduction in the computational cost compared to some existing related work. Finally, we demonstrate the effectiveness of our proposed algorithm through some numerical experiments, highlighting its efficiency and clear advantages.

Odirachukwunma Ugwu, Olaniyi Iyiola
Morgan State University
odirachukwunma.ugwu@morgan.edu,
olaniyi.iyiola@morgan.edu

## PP1

**Implementation of the Parareal Algorithm in the Open-Source Nektar++ Spectral/hp Element Framework**

Decomposing the spatial domain into multiple subdomains, each to be solved concurrently on individual processors, has become the traditional approach to reduce computational time for the numerical simulations of complex and large problems in high-performance computing. However, due to communication overhead, the maximum possible computational speed-up reaches an upper limit when the problem size per processor has become too small for distributed memory supercomputer architecture. In this respect, parallel-in-time algorithms are increasingly recognized as a promising solution to increase computational concurrency after the maximum speed-up has been achieved from pure spatial parallelism. One popular time-parallel approach is the Parareal algorithm of Lions et al. (2001) which exploits a fine and coarse time integrators in combination with an iterative procedure to achieve parallelism in time. Resulting from its simplicity, flexible and non-intrusive nature, as well as its applicability to both linear and non-linear problems, the Parareal algorithm has become one of the most popular parallel-in-time integration methods found in the literature. The implementation of the Parareal algorithm in the open-source Nektar++ spectral/hp element framework is described in this work. The extension of the MPI topology to allow concurrency in time and the implementation of a new driver subclass are presented. Applications of the Parareal algorithm to various problems are demonstrated.

Jacques Y. Xing
King's College London
jacques.xing@kcl.ac.uk

David Moxey
Department of Engineering
King's College London
david.moxey@kcl.ac.uk

Chris Cantwell
Imperial College London

c.cantwell@imperial.ac.uk

**PP1**
**Evaluation Graphics Processing Units (GPUs) Performance of Bisicles Ice-Sheet Flow Solver**

We evaluate the graphics processing unit (GPU) performance of the Berkeley-ISICLES (BISICLES) ice-sheet flow solver. Simulations that provide insights into future sea level rise require ensembles of simulations using high-resolution ice-sheet models, which can be computationally expensive. Solving time-independent stress balance equations to predict the ice velocity is the most computationally expensive part of ice-sheet simulations, both in terms of computer memory and execution time. To improve computational efficiency, we applied a Jacobian-free Newton-Krylov (JFNK) nonlinear solver using the Algebraic Multigrid (AMG) methods in the PETSc and Hypre libraries, running on AMD CPUs and NVIDIA GPUs. We focused on benchmark simulations for the evolution of continental scale Antarctica with three choices of finest spatial resolution near the grounding line; 2 km, 1 km, and 500 m. We demonstrate that GPU hardware capabilities can alleviate the computational cost and enable higher simulation throughput for future sea level rise predictions.

Zhengrong Zou, Anjali Sandip
University of North Dakota
zhengrong.zou@und.edu, anjali.sandip@und.edu

Daniel Martin
Lawrence Berkeley National Laboratory
dfmartin@lbl.gov

Hans Johansen
Lawrence Berkeley National Laboratory
Computational Research Division
hjohansen@lbl.gov