

Flash Talks – Session 2

Energy-Efficient Computing for Science Workshop

September 9-12, 2024
Bethesda, MD



U.S. DEPARTMENT OF
ENERGY

Office of
Science

Brian Austin, Berkeley Lab / NERSC

Orchestrating Power on Energy-Efficient Supercomputers

Opportunity and Proposed Research Direction

- What do you propose to do?
 - Manage power use within the supercomputer. Steer power toward jobs that need it and away from jobs that don't.
- Why now?
 - Without power management, the power-delivery for HPC systems is provisioned for maximum (not actual) power use. The cost of the resulting overinvestment in facility infrastructure is growing rapidly.

Execution and Timeline

- What are some of the key steps along the way?
 - Understand applications' power use and sensitivity to power reduction.
 - Define criteria to allocate power among jobs.
 - Develop robust and scalable software to implement these policies.
- What barriers to success can you anticipate?
 - Resistance to deploying power-management policies with uncertain performance consequences

State of the Art and Challenges

- How is this done today, what are the shortcomings?
 - Technologies to measure and manage power at the node level are widely available.
 - Tools for managing power across many nodes are not mature.
- Why is overcoming these challenges difficult?
 - Resource management (i.e. job scheduling) is a complex problem and is only complicated by the additional dimension of power.

Potential Impact

- What difference would a breakthrough in this area make?
 - Forestall a possible scenario where future supercomputing capabilities were limited by the availability of power and cooling rather than the advancement of computing technologies.
 - Increased computational throughput on a fixed facility budgets.
- How will we measure success?
 - Production HPC systems operate at or below prescribed power levels with minimal deviation from the target power.

Jean Luca Bez, Lawrence Berkeley Laboratory

Data Movement Orchestration for Energy-efficient Workflow Execution

Opportunity and Proposed Research Direction

- Scientific workflow co-design must focus on data movement as much as compute efficiency thus we propose three areas of focus towards data movement co-design:
 - I. Introspectable Programming Interfaces for productivity, performance, and energy efficiency
 - II. A data movement runtime system for energy-efficient orchestration
 - III. Novel data reduction and feature extraction strategies

State of the Art and Challenges

- With a focus on an IRI, data movement in wide-area networks presents opportunities for improving energy efficiency:
 - I. How to balance energy efficiency, performance, and programmer productivity without burdening workflow/application developers?
 - II. How to optimize the energy consumption of workflows in a computing environment where data and specialized computation are located in different geographic locations?
 - III. How to reduce the flow of low-value data?

Execution and Timeline

- AI-enabled intent-based programming for data/computation movement scheduling and querying:
 - I. Monitoring and estimating energy usage of the different system components
 - II. Runtime decisions to select the energy-efficient routes for data/computation movement and placement
 - III. Proactively extract vital information (small in size), helping reduce data movement between workflow stages

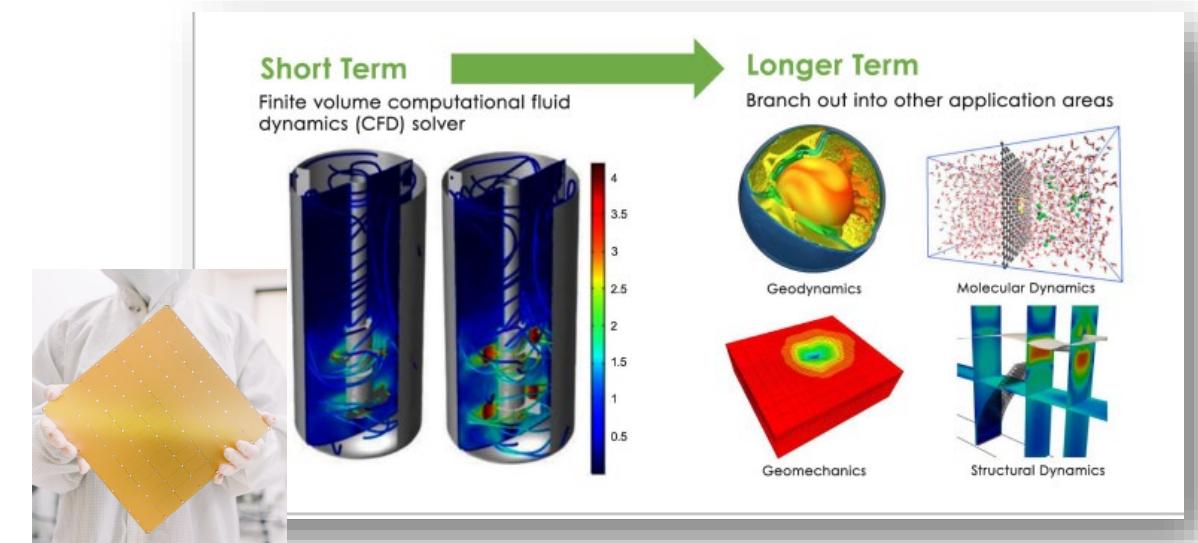
Potential Impact

- Besides improving energy efficiency while moving/accessing data that is geographically distant, a co-design approach would:
 - Enable faster scientific discovery with distributed/integrated data
 - Balance energy, performance, and productivity in data movement
 - Enabling intelligent and proactive data movement optimizations

Data Movement Orchestration for Energy-efficient Workflow Execution
Suren Byna (OSU, LBNL) Jean Luca Bez (LBNL), Houjun Tang (LBNL), and Wei Zhang (LBNL)

Opportunity and Potential Impact

- What: WSE-enabled scientific computing capability
- Why Now: WSE-3 hardware released in March 2024
- Impact: 100-1000x power efficiency, time-to-solution
- Metrics: Unstructured grid modeling; strong & weak scaling benchmarks



State of the Art and Challenges

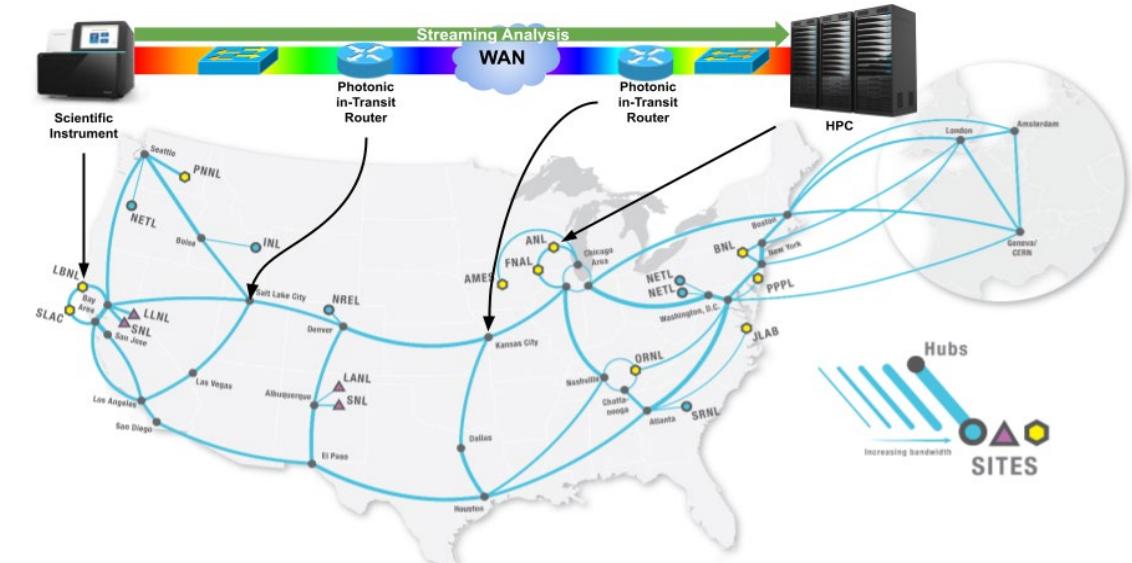
- State of art: On-chip processing power doubles every 2-2.5 years; power efficiency has not kept pace (inter-device bandwidths and latencies)
- Challenge: Paradigm shift to near-memory computing

Execution and Timeline

- Key steps: Components to solve linear system of equations; inter-wafer communications library
- 10-Year Roadmap: CS building blocks for distributed memory processing; coupled AI-physics for more gain
- Barriers: Cost, paradigm shift, high level programming
- Resources: public-private testbed; larger set of domain challenges

Opportunity and Potential Impact

- Can we offload computation to photonic in- and near-network devices for energy efficient computation?
- Photonic integrated chips (PICs) and in-network computing as enabler for this vision
 - A future in which photonic devices in-transit actively participate in energy efficient computation for Science
- Can an entire data processing pipeline to be photonics based?



ESnet map source: <https://www.es.net/news-and-publications/welcome-esnet6/esnet6-maps/>

State of the Art and Challenges

- Current practice is to move data from instruments to HPC facilities that consume Megawatts of power
 - It has been shown that under high message rates, in-network computing quickly becomes more power efficient than host-based solutions [1]
- New Science workloads will generate data at higher rates and at higher volumes, thus they will require more compute and may consume more energy
 - It has been recently demonstrated that a large-scale photonic chiplet can perform 160 TOPS/W for AI workloads [2]

Execution and Timeline

- Develop photonic technologies for both computation and interconnects (ex. optical packet switching)
- Rigorous metrics for in-network computing energy efficiency (5 years), in-network photonic computing (10 years), scientific workflow leveraging photonic in-transit computing (15 years)
- Co-design activities to integrate the hardware advances in PIC with the systems and software advances on in-network computing

[1] Y. Tokusashi, H. T. Dang, F. Pedone, R. Soulé, and N. Zilberman, "The case for in-network computing on demand," in Proceedings of the Fourteenth EuroSys Conference 2019, EuroSys '19, (New York, NY, USA), Association for Computing Machinery, 2019.

[2] Z. Xu, T. Zhou, M. Ma, C. Deng, Q. Dai, and L. Fang, "Large-scale photonic chiplet taichi empowers 160-tops/w artificial general intelligence," *Science*, vol. 384, no. 6692, pp. 202-209, 2024.

Ben Feinberg, Sandia National Laboratories

Modeling and Simulation for Hybrid Analog-Digital Systems

Opportunity and Proposed Research Direction

- Novel accelerator technologies have the potential to provide multiple order of magnitude improvements over conventional digital systems. But...
- Most non-trivial applications have some operations that are best handled digitally.
- We need simulators to enable research on combining novel accelerators with digital systems.

State of the Art and Challenges

- Current novel accelerator proposals often use hand-tuned applications and/or highly custom software stacks.
- These approaches suffer from:
 - Poor scalability to large applications
 - Limited usability by other researchers
 - Lots of reinventing the wheel

Execution and Timeline

1. Simulators that can run integrate programmable digital logic (e.g. CPU) with novel accelerator blocks.
2. With simulators we can start developing compilers, runtimes, and programming models to enable non-experts to use these systems.
3. With a software stack application developers can begin exploring how to best implement their algorithms.

Potential Impact

- Integration of novel accelerators and digital systems are essential for the success of these accelerators.
- Simulation frameworks can be common ground for hardware, systems, and applications researchers enabling true co-design.

Joseph S. Friedman, The University of Texas at Dallas

Non-Volatility for Logic: Do the Energy Benefits Outweigh the Energy Costs?

Opportunity and Potential Impact

- Emerging nanotechnology devices have the potential to revolutionize computing by decreasing energy consumption.
- Non-volatility is intriguing due to the suggested potential to save energy at circuit and system level.
- **But does this outweigh the device switching energy cost?**

State of the Art and Challenges

- Conventional computing architectures do not leverage non-volatility for energy-efficiency.
- A new paradigm for general purpose computing may eventually be invented that leverages non-volatility for energy efficiency.
- **Non-volatility is generally detrimental to switching energy.**

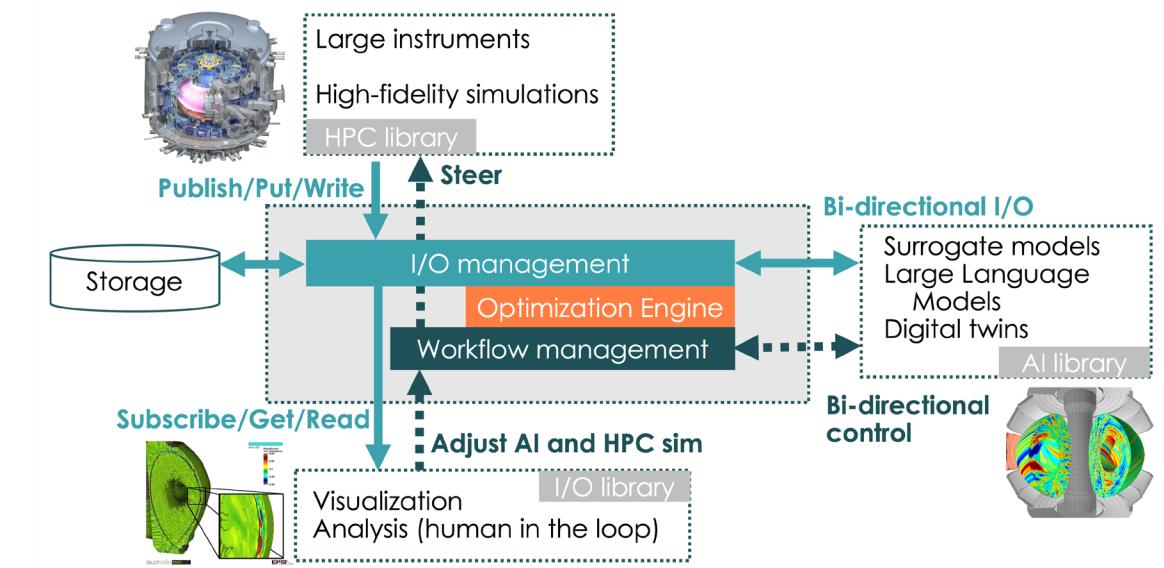


Execution and Timeline

- Non-volatility should be used for neuromorphic computing, programmability, and intermittent power
- **Non-volatility should be considered as a detrimental factor when evaluating novel device technologies for logical computing**
- New device technologies, both volatile and non-volatile, should be explored for logical computing and holistically evaluated

Opportunity and Potential Impact

- Deploy next generation workflows at scale
 - Performant, energy efficient and with resource constraints
 - Integrating AI
 - Allowing visualization and analysis on the fly
 - Handling PB data per second
 - Adapting to the profile and needs of each application



State of the Art and Challenges

- Current solutions have focused on increasing the FLOPS per Watt for individual HPC applications
 - Modern HPC apps are evolving into complex workflows with AI increasingly embedded
 - Complex queries are required to balance storage/compute
- Data access and resource allocation are multi-objective optimizations across multiple layers

Execution and Timeline

- Going beyond single applications. There is a need of
 - Re-designed algorithms that work on data knowledge instead of raw data and focus on interoperability
 - Adaptable workflows that can monitor the energy profile of applications and adjust resources and data transfers
 - Re-design querying capabilities for the needs of both traditional analysis/viz and new AI technologies

Paul Hovland, J. Hückelheim, S.H.K. Narayanan, Argonne National Laboratory

Energy-Efficient Derivatives, Derivatives for Energy Efficiency

Opportunity and Proposed Research Direction

- There exists an opportunity to refocus automatic differentiation (autodiff) research on energy efficiency.
- Simultaneously, derivatives provide insight into the parts of a computation that the final result is most sensitive to.
- Timeliness: autodiff and differentiable programming are readily available through frameworks such as PyTorch and JAX; energy consumption is a major concern

State of the Art and Challenges

- Current strategies for computing derivatives often seek to minimize flops, as a surrogate for minimizing time
- Future strategies may need to consider data movement as well in order to reduce energy consumption
- Balancing energy savings against increases in time or possible reductions in accuracy is a major challenge

Execution and Timeline

- New checkpointing/rematerialization strategies for trading recomputation against storage are needed
- Methods to compress checkpoints and “pre-accumulated” partial derivatives are needed
- Extensions to methods such as ADAPT and HAWQ that guide the introduction of reduced precision or approximate computations are needed

Potential Impact

- Reverse mode autodiff (backprop) is the workhorse of DNN training; any reduction in the energy consumption of DNN training could have a massive impact
- Autodiff is also mainstream for simulations and SciML
- On the right track if we achieve similar levels of accuracy in computations at reduced energy budgets

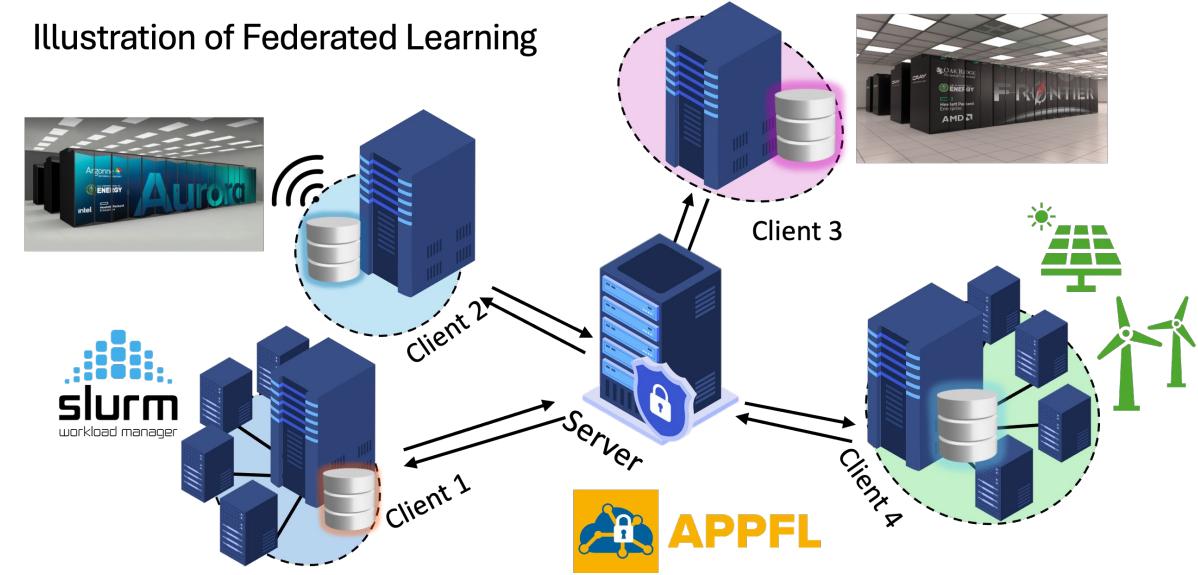
Kibaek Kim, Argonne National Laboratory

Energy-Efficient Computing in Federated Learning

Opportunity and Potential Impact

- The current energy demands in ML and FL are unsustainable.
- Needs for adaptive energy-aware algorithms and communication techniques.
- Unique opportunity to leverage renewable energy to address energy inefficiencies, aligning with global sustainability goals.
- Potential improvement in energy-to-performance ratios and lower operational costs for large-scale FL systems.
- Potential integration of renewable energy with effective communication protocols and flexible training models.

Illustration of Federated Learning



State of the Art and Challenges

- FL largely operates on homogeneous edge devices with limited focus on energy efficiency.
- Existing techniques like client sampling are insufficient for heterogeneous, cross-silo settings.
- High energy consumption in both computation and data transfer, with energy peaks stressing grid infrastructure.
- Overcoming challenges is difficult due to heterogeneous computing resources and the complex balance between performance and energy efficiency.

Execution and Timeline

- R&D of the capabilities to query various energy footprints from computing tasks.
- R&D of the energy-aware algorithms that dynamically schedule FL tasks based on energy efficiency metrics.
- Integration of FL systems with renewable energy sources (via signals from grid) to enhance sustainability.
- Complexity in modeling energy-performance trade-offs and the need for large-scale testbeds with real-world data.

Energy Efficient Probabilistic Scientific Computing on Neuromorphic Hardware



J. Darby Smith, Michael J. Schmidt, William Severa, **Hemanth Kolla** (Sandia National Labs)

Probabilistic approaches on neuromorphic hardware will accelerate scientific simulations.

- Develop probabilistic algorithms (coupled PDEs + SDEs) for scientific applications (i.e., go beyond deterministic PDEs).
- Leverage *neuromorphic advantage* for energy efficient random walks.
- Comparable accuracy but orders-of-magnitude more efficient for **post-exascale** scientific computing.
- Assess accuracy-cost tradeoff for benchmark problems.
 - Identify R&D for increased scale & physics complexity.

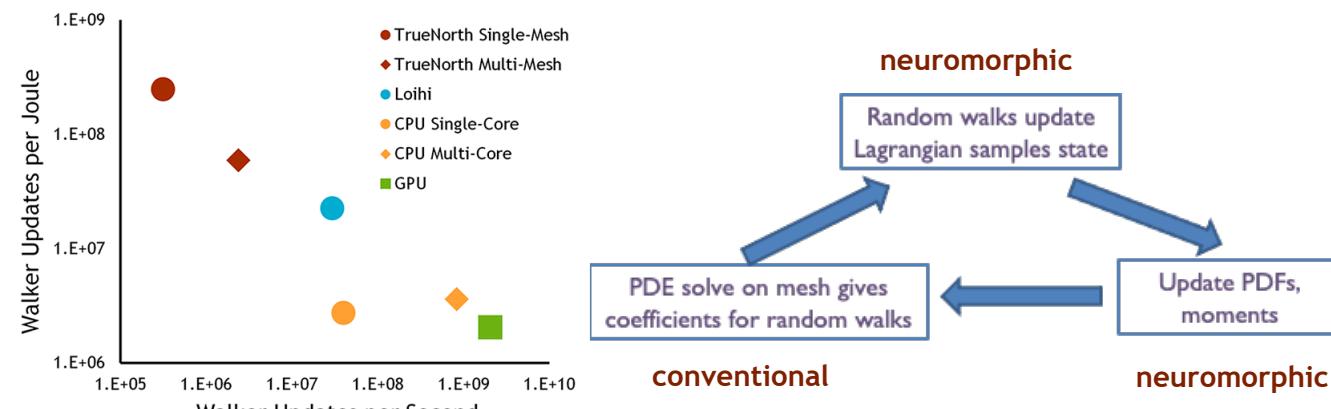


Image: Smith et al., Nature Electronics 2022.

State of the Art and Challenges

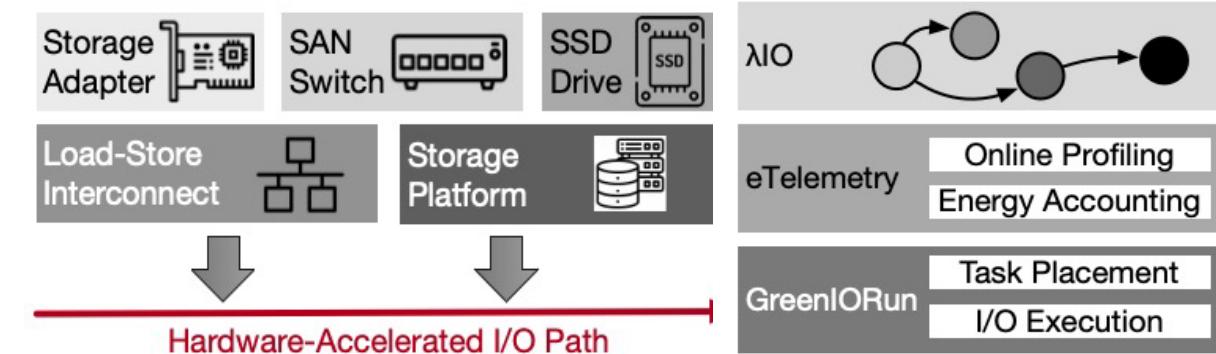
- Classical approaches for multi-scale problems (e.g., turbulence) adopt deterministic Eulerian PDE discretizations.
- Grid/timestep requirements scale exponentially.
- Probabilistic approaches (e.g., solving PDF equations) use ensemble random walks \Rightarrow not efficient on conventional.
- Heterogeneous **hardware + algorithm co-design** is necessary.

Execution and Timeline

- Exemplar multi-scale problem:
 - Turbulence simulations with “transported PDF” approach (beyond DNS/LES).
- Proof-of-concept on 1st-gen neuromorphic (2 yrs), specializing hardware/algorithms for next-gen (5 yrs), **massively parallel simulations** (10-15 yrs).
- Neuromorphic testbeds and ecosystem at **Leadership Computing Facilities** will be critical.

Opportunity and Potential Impact

- Hardware-accelerated I/O from the NVMe drive and storage adapter to switching fabric and interconnects
- A revamped storage stack that can automatically refactor the I/O processing over on-path accelerators based on energy appraisal
- Orders of magnitude higher IOPS/Joule



State of the Art and Challenges

- Storage servers have become dense and power-hungry
- Existing storage stacks apply a decade-old “smart-sender dumb-receiver” design philosophy
- Require dozens of cores to busy-drive the I/O parallelism and fully use the storage bandwidth, yielding dramatic power consumption (>500W)

Execution and Timeline

- Key steps: (a) perform energy accounting at the per-IO granularity; (b) reconstruct the I/O data path on the fly and map tasks to a suitable hardware substrate; (c) orchestrate the I/O execution without exceeding the computing limit of the corresponding devices
- Incrementally integrate emerging I/O accelerators and transparently expose them to the software storage layer

Matt Menickelly, Argonne National Laboratory

Energy-Efficient Algorithms for Nonlinear Optimization

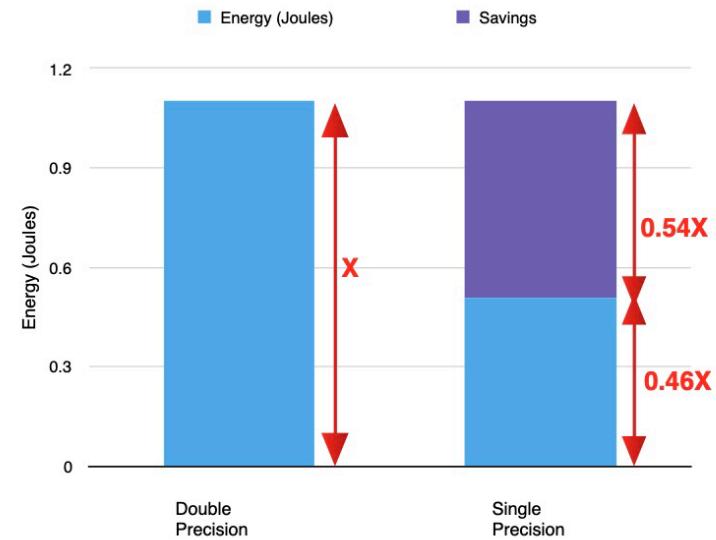
Opportunity and Potential Impact

- Nonlinear optimization (NLO) problems are fundamental to many critical DOE applications (e.g., design and control of power-grid networks, design of experiments, inversion of material properties from light-source data)
- Given energy consumption concerns, we must develop algorithms for NLO that maintain nearly deterministic guarantees, but substantially reduce energy-to-solution

State of the Art and Challenges

- SOTA algorithms are designed to optimize compute time, iteration complexity, memory footprint; do these directly translate to energy savings? Bottleneck of virtually any NLO method is the iterative solution of linear systems.
- Research is ongoing in how to employ 1) mixed precision, 2) randomization (e.g., sketching) and 3) asynchronicity to solve these systems. Better hardware codesign will be essential in implementing these effectively.

Figure: From “Doing Moore with Less ...” by Leyffer et al, 2016. A simple implementation of a single-double precision switching NLO algorithm (right bar) generates energy savings of 54% on a reasonable benchmark.



Execution and Timeline

- Orchestrate *hardware design* to facilitate mixed precision/energy consumption estimates/structured random matrices/asynchronous communication, *software design* to expose and make high-level routines to employ these hardware features, and *fundamental algorithmic research*
- Prototype stacks could take ~5 years of development

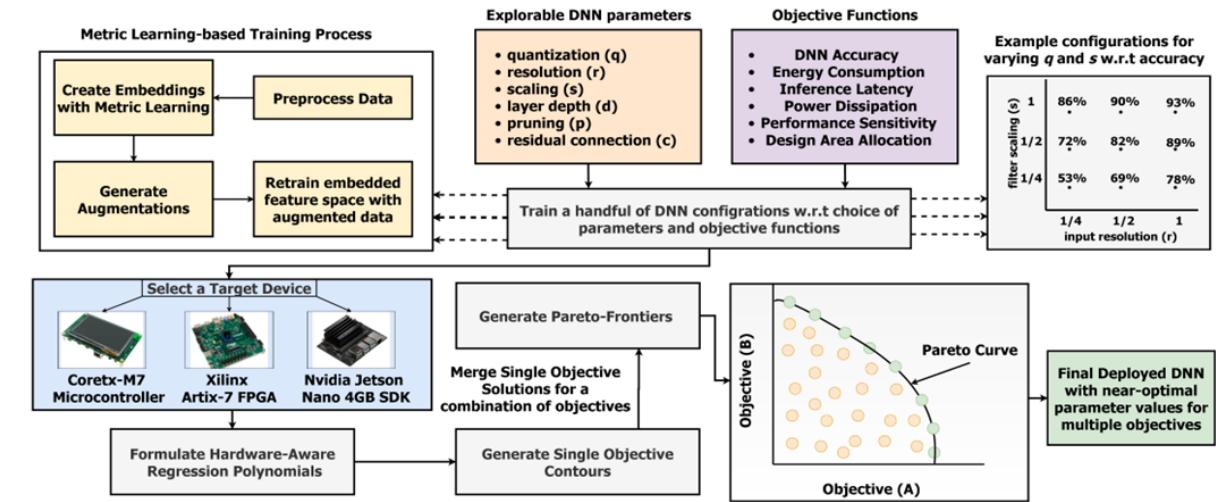
Enhancing Energy Efficiency: Addressing Energy Costs of Communication through Hardware-Software Codesign

Opportunity and Potential Impact

- Proposed Regression-based search enables efficient and fast architecture selection for energy-optimized AI models by fine-tuning parameters such as filter scaling, quantization, and pruning.
- Provides efficient co-design at edge in various instrumentation across metrology, sensing, synthesis and detection.
- Communication-related energy costs can be minimized by storing synaptic weights and broadcasting only non-zero partial sums.
- This approach enhances deployment flexibility, allowing users to explore performance profiles and select near-optimal configurations suited to hardware specific requirements.

State of the Art and Challenges

- Traditionally, this fine-tuning has been done through methods like grid search, random search, or more advanced neural architecture search (NAS) techniques, which can be resource-intensive and time-consuming.
- Current search methods heavily rely on GPU usage, potentially hindering the development of AI hardware, software, and applications.
- The extreme resource utilization bottleneck directly affects the potential for AI-enabled scientific discovery within the Department of Energy (DOE) and beyond.

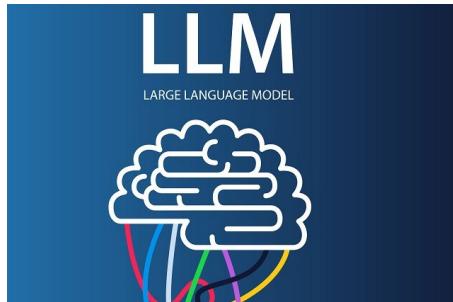


Execution and Timeline

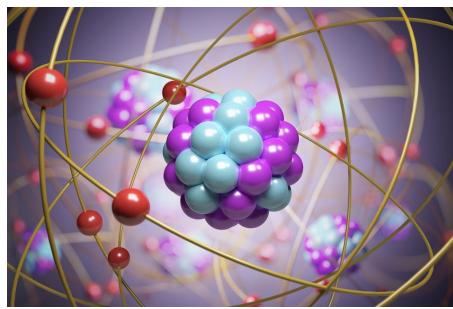
- Provide analytical and fast architecture search for large and complex models on variety of hardware, data and application domains.
- The proposed solution could evolve to an autonomous architecture selection process, enabling real-time scalability and deployment while fostering collaboration between AI, high-performance computing and quantum technologies for scientific and industrial advancements.
- This research over the next 15 years will foster interdisciplinary collaboration, provide access to cutting-edge hardware and diverse datasets, and address barriers such as hardware constraints and scalability.

Doru Thom Popovici, Lawrence Berkeley National Lab

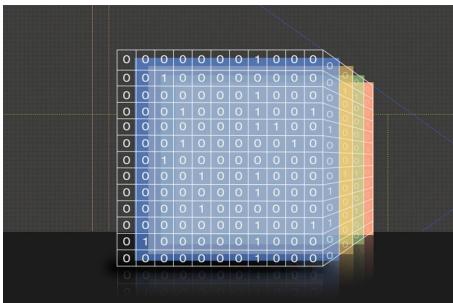
Reshaping Algorithms and Data Movement for Energy Efficient Hardware Acceleration



AI Algorithms



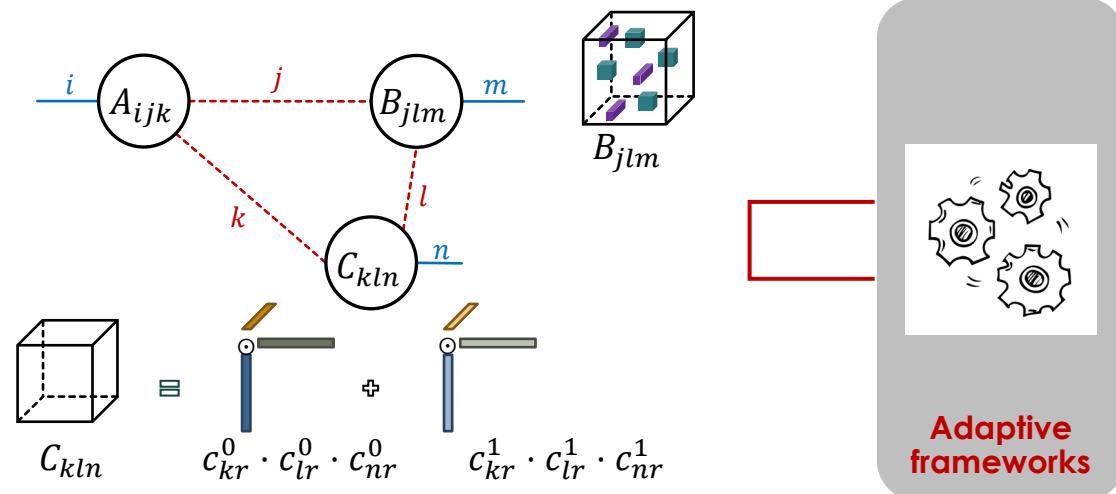
Physics and Chemistry Simulations



Data Science

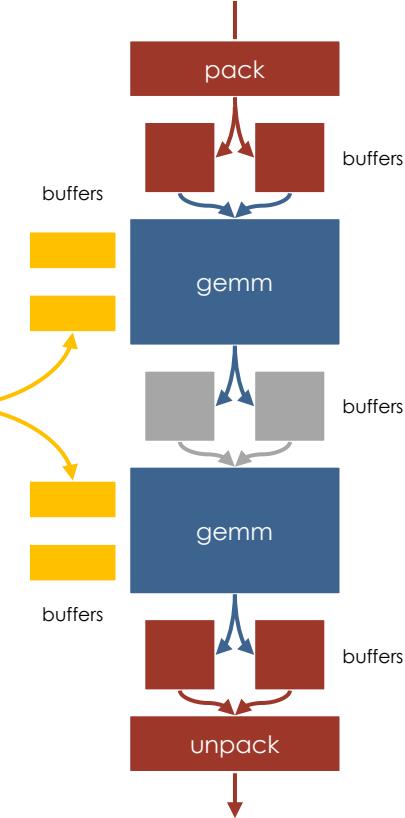
Problem

- A. Data movement accounts for a **large percentage of power consumption**.
- B. Scientific simulations and ML algorithms need **ever increasing amounts of data**.



Approach

- A. **New algorithms and mathematical models** to enable more aggressive optimizations.
- B. Hardware/software solutions for **data reduction when communication** is required.
- C. **Data driven approaches** to adaptively optimize computation and data movement.



Collaborators: Roel Van Beeumen, Angelos Ioannou, Mauro Del Ben, Mario Vega, Meriam Gay Bautista-Jurney, Fabien Chaix, Xiaokun Yang, John Shalf



BERKELEY LAB
Bringing Science Solutions to the World

Matthew D. Sinclair, UW-Madison (co-Authors: Bobby Bruce, William Godoy, Oscar Hernandez, Jason Lowe-Power, and Shivaram Venkataraman)

Creating Flexible, High Fidelity Energy Modeling for Future HPC Systems

Opportunity and Potential Impact

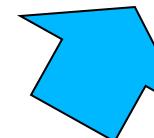
- **Energy efficient, large scale, co-designed HPC sys required**
- Typically sim & modeling tools enable early-stage design exploration ... but energy modeling SOTA lagging
- **Key Question:** How to Design energy models to scalably represent modern systems, pre-tapeout, with high fidelity?
- **Proposal:** develop credible, open-source, high fidelity energy models for current & future systems

State of the Art and Challenges

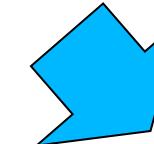
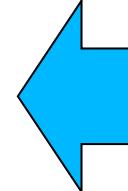
- **Early-stage energy tools divided, arch-specific, out-of-date**
 1. First-principal tools (e.g., McPAT) 8+ years old
 2. Empirical measurements hard to generalize
 3. ML-based models suffer extrapolating to new HW
 4. Tools based on tapeout values time consuming, expensive, only late in design process
 5. Low-level Spice models accurate, but require proprietary information, hard to scale to large systems

Wither Moore's Law?

Faster hardware



Larger Datasets



Improved algorithms
(e.g., deeper DNNs)

Moore's Law: **virtuous cycle** of progress in many fields
... slowing of Moore's Law **threatens progress**
Modern apps have exponential compute, power needs

Execution and Timeline

- Proposed Steps:
 - **Current Sys:** Fine-grained info → domain scientists usable
 - **Future Sys:** leverage burgeoning open-source HW + learn and predict patterns
 - **Ease of Use:** integrate with popular tools like gem5+SST
- **Long-Term Vision:** model energy model as easy as perf. model
- Enable energy-aware co-design to optimize HPC systems **early**
- **Developing resources will take time, DOE can spur it**



Opportunity and Potential Impact

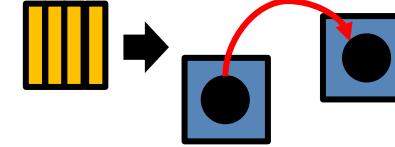
- **Software:** Compute where data is, **Circuits:** 3D integration with heterogeneous technology integration, **Architecture:** CXL + near-data processing with support for sparse computation, **Devices:** reconfigurable transistors
- Addresses key limitations of current technologies (excessive data movements, device scaling, irregularity)

State of the Art and Challenges

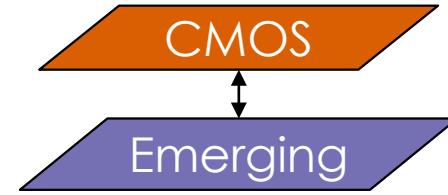
- High energy cost of moving data to compute, most hardware designed for regular workloads (GEMM), CMOS scaling hitting limitations
- Traditional approaches, limited applicability to a single layer-**ineffective**; crucial to break through the boundaries of software, hardware, devices, and materials to innovate further

Software: Intelligent Runtimes

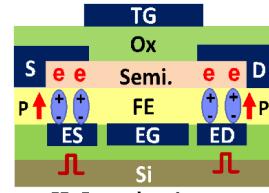
Send Compute to Data



Circuits: 3D Integration



Devices: Emerging Material, Device Technologies



Hardware: CXL, NDP, Sparse computation



Execution and Timeline

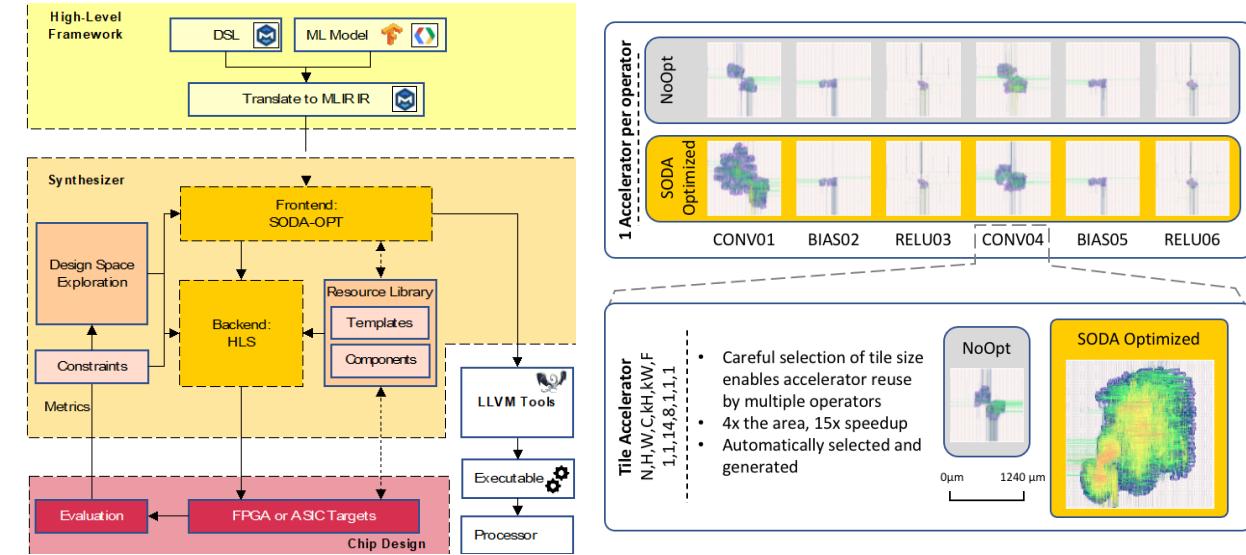
- (1) Define key domain combinations, (2) define workload characteristics, (3) establish cross-domain expert group, (4) define targets from each technology layer, and (5) get crackin!
- Years 1-10: ready software, architecture, POC devices/materials; Years 11-15: large-scale production-ready technology stack

Antonino Tumeo, PNNL

Agile Hardware Synthesis for Energy Efficient Systems

Opportunity and Potential Impact

- Domain-specific systems are the primary way to keep improving energy efficiency
- Developing custom accelerators by hand is complex, and feasible only for applications with mass market appeal
- Open-source, modular, extensible, compiler-based hardware synthesis tools from high-level programming frameworks to silicon can bridge the productivity gap
- Chiplets enable development of composable and scalable custom systems that can address DOE's needs such as autonomous experimental workflows that performs data analytics, artificial intelligence, and simulation across a continuum of computing



State of the Art and Challenges

- The co-design process has been successful, but can only be applied for common computing patterns and applications with broad market adoption due to its complexity
- Several approaches to hardware synthesis/generation available, but they do not cover the entire design stack and/or are openly accessible

Execution and Timeline

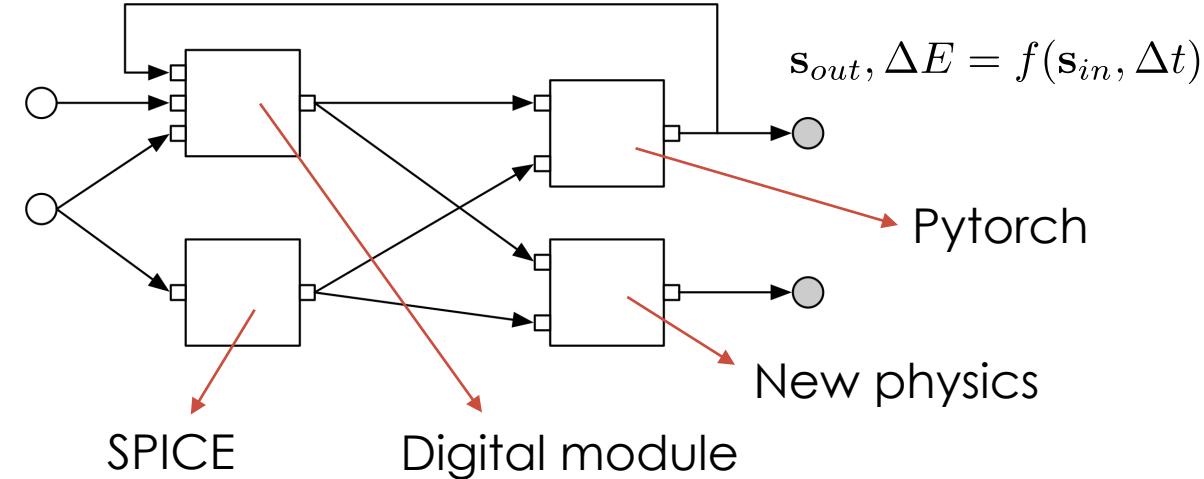
- Creation of an open hardware technology common
- Develop new synthesis and design space exploration methodologies that can enable optimization across multiple metrics (power, performance, area, and other)
- Co-existence and co-operation with proprietary tools and intellectual properties is critical
- Access to data open and proprietary tools and prototyping facilities are critical to develop and evaluate new methodologies and designs

Opportunity and Potential Impact

- We have 20+ years of research on emergent devices and small prototype architectures to build on.
- Developing new tools capable of assessing the potential of emergent devices and novel compute approaches to solve problems at scale will help us identify new pathways to energy efficient computing.
- It will help establish clear target for novel hardware

State of the Art and Challenges

- We don't have simple ways of estimating the power consumption of heterogeneous architectures interfacing digital designs with emergent technologies.
- There are significant challenges: power distribution, ADCs, interface to/from digital, need to bridge dissimilar timelines, packaging! (chiplet, tiling) analyze performance of complex workflows and establish comparison with digital baseline implementations.



Through Threadwork (DOE microelectronics) we explored novel approaches to model heterogeneous architectures in a massively parallel way. We need better, more scalable solutions.

Execution and Timeline

- We need better tools to integrate heterogeneity in hardware design
- We need consistent methods to compute energy consumption of digital designs
- We need 21st century simulation pipelines
- We need to focus on interfaces with digital

Accelerating Energy-Efficient Scientific High-Performance Computing with Hardware Specialization

Opportunity and Proposed Research Direction

- **What?** Develop energy-efficient, specialized hardware accelerators for HPC workloads leveraging agile prototyping, true co-design and chiplet technology.
- **Why Now?** Moore's Law slowdown, advancements in chiplet technology, and the rise of open hardware standards (RISC-V) create timely opportunities.

State of the Art and Challenges

- **Current Approaches:** the current HPC uses CPUs and GPGPUs, which aren't energy-efficient for all tasks.
- **Challenges:** Lack of agile hardware prototyping, expertise in accelerator development, and complex verification processes make it hard to adopt specialized hardware. Overcoming these challenges requires collaboration and innovation.

Execution and Timeline

- **Key Steps:** Develop agile hardware prototyping methodologies, integrate chiplets for energy efficiency, bridge HPC-hardware expertise gap.
- **Timeline:** 5 years: Establish methodologies. 10 years: Integrate accelerators into HPC systems.
- **Barriers:** Complexity in chip development such as verification, a longer time to market, expertise gaps.

Potential Impact

- **Breakthrough Impact:** Significant reductions in energy consumption by several order of magnitude and faster, more efficient HPC systems.
- **Measuring Success:** Progress will be measured by energy efficiency gains, successful hardware integration, and shorter prototyping cycles.