# Preparing for Post-Exascale Computing

**Jeffrey S. Vetter**

*With many contributions from ACSR Section and Colleagues*

The 30th Anniversary Symposium of the Center for Computational Sciences
University of Tsukuba
13 Oct 2022

*Congratulations on 30th Anniversary!*

ORNL is managed by UT-Battelle, LLC for the US Department of Energy

**U.S. DEPARTMENT OF ENERGY**

https://www.ornl.gov/section/advanced-computing-systems-research (https://j.mp/acsrs)
vetter@computer.org

# Highlights

- 15 years to go from Exascale ideation to deployment
  - Reports and predictions

- Current status
  - Systems (Frontier)
  - Exascale Computing Project

- Exascale: What did we get right, get wrong, overlook?

- Post Exascale?
  - Heterogeneous systems enabled by Heterogeneous integration and Chiplets
  - Codesign becomes even more important

- Abikso: Microelectronics codesign project

# Exascale Reports (and predictions) from 2007 to 2014



**Modeling and Simulation at the Exascale for Energy and the Environment**

Report on the Advanced
Town H...
Simulation and Modeling at the Exa...
and Glob...

Co-Chairs: Lawrence Berkeley Nation...
Oak Ridge National Labor...
Argonne National Laborat...

Office of
Advanced
Scientific Computing
Research Contact: Michael Strayer

Special Assistance
Technical: Lawrence Berkeley Nation...
Deb Agarwal, David Bailey,
William Collins, Nikos Kyrpi...
Peter Nugent, Leonid Olike...
Lin-Wang Wang, Michael W...

Oak Ridge National Labor...
Eduardo D'Azevedo, David...
James Hack, Victor Hazlew...
Bronson Messer, Anthony M...
B. (Rad) Radhakrishnan, N...
Jeffrey Vetter, Gilbert Weig...

Argonne National Laborat...
Raymond Bair, Pete Beckm...
Ed Frank, Ian Foster, Willia...
Robert Jacob, Kenneth Ker...
Jorge Moré, Lois McInnes,...
Michael Papka, Robert Ros...

Administrative: Lawrence Berkeley Nation...
Oak Ridge National Labor...
Argonne National Laborat...

Publication: Oak Ridge National Labor...
Argonne National Laborat...

Editorial: Oak Ridge National Labor...
Argonne National Laborat...

This report is available on the web at http://www.sc.c...

---

**ExaScale Computing Study: Technology Challenges in Achieving Exascale Systems**

Peter Kogge, Editor & Study Lead
Keren Bergman
Shekhar Borkar
Dan Campbell
William Carlson
William Dally
Monty Denneau
Paul Franzon
William Harrod
Kerry Hill
Jon Hiller
Sherman Karp
Stephen Keckler
Dean Klein
Robert Lucas
Mark Richards
Al Scarpelli
Steven Scott
Allan Snavely
Thomas Sterling
R. Stanley Williams
Katherine Yelick

September 28, 2008

This work was sponsored by DARPA IPTO i...
as Program Manager; AFRL contract number...
interest of scientific and technical informatio...
Government's approval or disapproval of its i...

Using Government drawings, specifications, ...
purpose other than Government procurement ...
The fact that the Government formulated or s...
does not license the holder or any other perso...
manufacture, use, or sell any patented inventi...

APPROVED FOR PUBLIC RELEASE, DIST...

---

**SCIENTIFIC GRAND CHALLENGES: CROSSCUTTING TECHNOLOGIES FOR COMPUTING AT THE EXASCALE**

**Report from the Workshop Held February 2-4, 2010**

Sponsored by the U.S. Department of Energy, O...
Research, Office of Science; and the Office of A...
National Nuclear Security Administration

*Chair*, **David L. Brown**
Lawrence Livermore National Laboratory

*Chair*, **Paul Messina**
Argonne National Laboratory

**Theme I: Domain Science and System Architecture**

*Principal Lead*, **David Keyes**
King Abdullah University of Science and Technology a...

*Co-Lead*, **John Morrison**
Los Alamos National Laboratory

*Co-Lead*, **Robert Lucas**
University of Southern California

*Co-Lead*, **John Shalf**
Lawrence Berkeley National Laboratory

**Theme II: System Software**

*Principal Lead*, **Pete Beckman**
Argonne National Laboratory

*Co-Lead*, **Ron Brightwell**
Sandia National Laboratories

*Co-Lead*, **Al Geist**
Oak Ridge National Laboratory

**Theme III: Programming Models and Environment**

*Principal Lead*, **Jeffrey Vetter**
Oak Ridge National Laboratory and Georgia Institute of Technology

---

**ASCAC Subcommittee for the Top Ten Exascale Research Challenges**

**Subcommittee Chair**
Robert Lucas (University of Southern California

**Subcommittee Members**
James Ang (Sandia National Laboratories)
Keren Bergman (Columbia University)
Shekhar Borkar (Intel)
William Carlson (Institute for Defense Analyses)
Laura Carrington (UC, San Diego)
George Chiu (IBM)
Robert Colwell (DARPA)
William Dally (NVIDIA)
Jack Dongarra (U. Tennessee)
Al Geist (ORNL)
Gary Grider (LANL)
Rud Haring (IBM)
Jeffrey Hittinger (LLNL)
Adolfy Hoisie (PNLL)
Dean Klein (Micron)
Peter Kogge (U. Notre Dame)
Richard Lethin (Reservoir Labs)
Vivek Sarkar (Rice U.)
Robert Schreiber (Hewlett Packard)
John Shalf (LBNL)
Thomas Sterling (Indiana U.)
Rick Stevens (ANL)

---

The International Journal of High
Performance Computing Applications
25(1) 3–60
© The Author(s) 2011
Reprints and permission:
sagepub.co.uk/journalsPermissions.nav
DOI: 10.1177/1094342010391989
hpc.sagepub.com
**$SAGE**

**The International Exascale Software Project roadmap**

**Jack Dongarra, Pete Beckman, Terry Moore, Patrick Aerts, Giovanni Aloisio, Jean-Claude Andre, David Barkai, Jean-Yves Berthou, Taisuke Boku, Bertrand Braunschweig, Franck Cappello, Barbara Chapman, Xuebin Chi, Alok Choudhary, Sudip Dosanjh, Thom Dunning, Sandro Fiore, Al Geist, Bill Gropp, Robert Harrison, Mark Hereld, Michael Heroux, Adolfy Hoisie, Koh Hotta, Zhong Jin, Yutaka Ishikawa, Fred Johnson, Sanjay Kale, Richard Kenway, David Keyes, Bill Kramer, Jesus Labarta, Alain Lichnewsky, Thomas Lippert, Bob Lucas, Barney Maccabe, Satoshi Matsuoka, Paul Messina, Peter Michielse, Bernd Mohr, Matthias S. Mueller, Wolfgang E. Nagel, Hiroshi Nakashima, Michael E Papka, Dan Reed, Mitsuhisa Sato, Ed Seidel, John Shalf, David Skinner, Marc Snir, Thomas Sterling, Rick Stevens, Fred Streitz, Bob Sugar, Shinji Sumimoto, William Tang, John Taylor, Rajeev Thakur, Anne Trefethen, Mateo Valero, Aad van der Steen, Jeffrey Vetter, Peg Williams, Robert Wisniewski and Kathy Yelick**

**Abstract**
Over the last 20 years, the open-source community has provided more and more software on which the world's high-performance computing systems depend for performance and productivity. The community has invested millions of dollars and years of effort to build key components. However, although the investments in these separate software elements have been tremendously valuable, a great deal of productivity has also been lost because of the lack of planning, coordination, and key integration of technologies necessary to make them work together smoothly and efficiently, both within individual petascale systems and between different systems. It seems clear that this completely uncoordinated development model will not provide the software needed to support the unprecedented parallelism required for peta/exascale computation on millions of cores, or the flexibility required to exploit new hardware models and features, such as transactional memory, speculative execution, and graphics processing units. This report describes the work of the community to prepare for the challenges of exascale computing, ultimately combing their efforts in a coordinated International Exascale Software Project.

**Keywords**
exascale computing, high-performance computing, software stack

**Table of Contents**

# My first trip to Tsukuba!

- **International Exascale Software Project**

- **Meeting 3:**
  **Tsukuba, Japan**
  **Oct. 18-20, 2009**

# DOE HPC Roadmap to Exascale Systems

| FY 2012 | FY 2016 | FY 2018 | FY 2021 | FY 2022 | FY 2023 |
|---------|---------|---------|---------|---------|---------|



decommissioned

**Titan**
**ORNL**
Cray/AMD/NVIDIA

**Mira**
**ANL**
IBM BG/Q

**Sequoia**
**LLNL**
IBM BG/Q

**Theta**
**ANL**
Cray/Intel KNL

**Cori**
**LBNL**
Cray/Intel Xeon/KNL

**Trinity**
**LANL/SNL**
Cray/Intel Xeon/KNL

**Summit**
**ORNL**
IBM/NVIDIA

**Sierra**
**LLNL**
IBM/NVIDIA

**FRONTIER**
**ORNL**
HPE/AMD

**Polaris**
**ANL**
HPE/AMD/NVIDIA

**Perlmutter**
**LBNL**
HPE/AMD/NVIDIA

**Aurora**
**ANL**
Intel/HPE

**CROSSROADS**
**LANL/SNL**
HPE/Intel

**EL CAPITAN**
**LLNL**
HPE/AMD

## Exascale Systems

ECP EXASCALE COMPUTING PROJECT

Version 2.0

# DOE HPC Roadmap to Exascale Systems



| FY 2012 | FY 2016 | FY 2018 | FY 2021 | FY 2022 | FY 2023 |
|---------|---------|---------|---------|---------|---------|

**Titan**
**ORNL**
Cray/AMD/NVIDIA

**Mira**
**ANL**
IBM BG/Q

**Theta**
**ANL**
Cray/Intel KNL

**Cori**
**LBNL**
Cray/Intel Xeon/KNL

**Summit**
**ORNL**
IBM/NVIDIA

To date, only NVIDIA GPUs

**FRONTIER**
**ORNL**
HPE/AMD

→ **Exascale Systems**

**Polaris**
**ANL**
HPE/AMD/NVIDIA

**Aurora**
**ANL**
Intel/HPE

**Perlmutter**
**LBNL**
HPE/AMD/NVIDIA

AMD, Intel and NVIDIA GPUs!

decommissioned

**Sequoia**
**LLNL**
IBM BG/Q

**Trinity**
**LANL/SNL**
Cray/Intel Xeon/KNL

**Sierra**
**LLNL**
IBM/NVIDIA

**CROSSROADS**
**LANL/SNL**
HPE/Intel

**EL CAPITAN**
**LLNL**
HPE/AMD

OAK RIDGE National Laboratory
Version 2.0

# Frontier System



**System**
- 74 compute racks
- 29 MW Power Consumption
- 9,408 nodes
- 9.2 PB memory
  (4.6 PB HBM, 4.6 PB DDR4)
- Cray Slingshot network with dragonfly topology
- 37 PB Node Local Storage
- 716 PB Center-wide storage
- 4000 ft² foot print

# Frontier Cabinet

**Olympus rack**
- 128 AMD nodes
- 8,000 lbs
- Supports 400 KW



# Frontier Node

**AMD extraordinary engineering**
- 1 AMD "Trento" CPU (optimized Milan)
- 4 AMD MI250X GPUs
- 512 GiB DDR4 memory on CPU
- 512 GiB HBM2e total per node
- 4 Cassini NICs connected to the 4 GPUs

**Compute blade**
- 2 AMD nodes



**All water cooled, even DIMMS and NICs**

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

# OAK RIDGE NATIONAL LABORATORY'S FRONTIER SUPERCOMPUTER

- **74 HPE Cray EX cabinets**

- **9,408 AMD EPYC CPUs, 37,632 AMD GPUs**

- **700 petabytes of storage capacity, peak write speeds of 5 terabytes per second using Cray Clusterstor Storage System**

- **90 miles of HPE Slingshot networking cables**

Sources: May 30, 2022 Top500 release

## TOP500
# #1

**1.1 exaflops** of performance on the May 2022 Top500.

## GREEN500
# #1, #2

62.04 gigaflops/watt power efficiency on a single cabinet.
52.23 gigaflops/watt power efficiency on the full system.

## HPL-AI
# #1

**6.88 exaflops** on the HPL-AI benchmark.

# Exascale Computing Project has three technical areas to meet national goals

**Performant mission and science applications @ scale**

| Foster application development | Ease of use | Diverse architectures | HPC leadership |

## Application Development (AD)
Develop and enhance the predictive capability of applications critical to the DOE

## Software Technology (ST)
Produce expanded and vertically integrated software stack to achieve full potential of exascale computing

## Hardware and Integration (HI)
Integrated delivery of ECP products on targeted systems at leading DOE computing facilities

25 applications ranging from national security, to energy, earth systems, economic security, materials, and data

80+ unique software products spanning programming models and run times, math libraries, data and visualization

6 vendors supported by PathForward focused on memory, node, connectivity advancements; deployment to facilities

https://www.exascaleproject.org/

EXASCALE
COMPUTING
PROJECT

# Application Development KPP-1 and KPP-2 Readiness Overview

| KPP-1 App | Aurora EAS (Intel Proprietary) | Frontier TDS |
|---|---|---|
| LatticeQCD | Verified | Improving Perf. |
| NWChemEx | Full Build/Test | Initial Build/Test |
| EXAALT | Verified | Improving Perf. |
| QMCPACK | Initial Build/Test | Improving Perf. |
| ExaSMR | Improving Perf. | Improving Perf. |
| WDMApp | Improving Perf. | Improving Perf. |
| WarpX | Verified | Improving Perf. |
| ExaSky | Improving Perf. | Improving Perf. |
| EQSIM | Initial Build/Test | Improving Perf. |
| E3SM-MMF | Improving Perf. | Improving Perf. |
| CANDLE | Ready | Improving Perf. |

| KPP-2 App | Aurora EAS (Intel Proprietary) | Frontier TDS |
|---|---|---|
| GAMESS | Improving Perf. | Improving Perf. |
| ExaAM | Initial Build/Test | Improving Perf. |
| ExaWind | Verified | Improving Perf. |
| Combustion-PELE | Initial Build/Test | Improving Perf. |
| MFIX-Exa | Verified | Improving Perf. |
| ExaStar | Full Build/Test | Improving Perf. |
| Subsurface | Stretch | Improving Perf. |
| ExaSGD | Stretch | Improving Perf. |
| ExaBiome | Stretch | Improving Perf. |
| ExaFEL | Full Build/Test | Blocked (ROCm) |

ExaFel (Blocked ROCm) – The project team noticed during a recent work around to compiler bugs preventing compilation of Spinifel relies on an upstream compiler built locally by the OLCF team. Integration into a formal ROCm release is still pending, so we have chosen to mark this as formally blocked even though the team can currently make progress with the unofficial compiler build.

# Extreme-scale Scientific Software Stack (E4S)

- <u>E4S</u>: HPC software ecosystem – a curated software portfolio

- A **Spack-based** distribution of software tested for interoperability and portability to multiple architectures

- Available from **source**, **containers**, **cloud, binary caches**

- Leverages and enhances SDK interoperability thrust

- Not a commercial product – an open resource for all

- Growing functionality: May 2022: E4S 22.05 – 100+ full release products

https://spack.io

Spack lead: Todd Gamblin (LLNL)

https://e4s.io

E4S lead: Sameer Shende (U Oregon)

Also includes other products, e.g.,
**AI:** PyTorch, TensorFlow, Horovod
**Co-Design:** AMReX, Cabana, MFEM

| | | |
|---|---|---|
| **Community policies** Commitment to SW quality | **DocPortal** Single portal to all E4S product info | **Portfolio testing** Especially leadership platforms |
| **Curated collection** The end of dependency hell | **Quarterly releases** Release 22.2 – February | **Build caches** 10X build time improvement |
| **Turnkey stack** A new user experience | https://e4s.io | **Post-ECP Strategy** Commercial E4S, SSO |

ECP EXASCALE COMPUTING PROJECT

# ECP is Improving the LLVM Compiler Ecosystem

| LLVM | + SOLLVE | + PROTEAS-TUNE | + FLANG | + HPCToolkit | + NNSA | Vendors |
|------|----------|----------------|---------|--------------|--------|---------|
| • Very popular open-source **compiler infrastructure**<br>• **Permissive** license<br>• **Modular**, well-defined IR allows use by a lot of different languages, ML frameworks, etc.<br>• **Backend infrastructure** allowing the efficient creation of backends for new (heterogeneous) hardware.<br>• A **state-of-the-art C++ frontend**, CUDA support, scalable LTO, sanitizers and other debugging capabilities, and more. | • Enhancing the implementation of OpenMP in LLVM<br><br>• Unified memory<br><br>• Prototype OMP features for LLVM<br><br>• OMP Optimizations<br><br>• OMP test suite<br><br>• Tracking OMP implementation quality<br><br>• Training | • Core optimization improvements to LLVM<br>  • OpenMP offload<br><br>• OpenACC capability for LLVM<br>  • Clacc<br>  • Flacc<br><br>• Autotuning for OpenACC and OpenMP in LLVM<br><br>• Integration with Tau performance tools<br><br>• SYCL characterizing and benchmarking<br><br>• Leading LLVM-DOE fork<br><br>• Training | • Developing an open-source, production Fortran frontend<br><br>• Upstream to LLVM public release<br><br>• Support for OpenMP and OpenACC<br><br>• Recently approved by LLVM<br><br>• Initial implementation of serial F77 compiler for CPUs under review | • Improvements to OpenMP profiling interface OMPT<br><br>• OMPT specification improvements<br><br>• Refine HPCT for OMPT improvements | • Enhancing LLVM to optimize template expansion for FlexCSI, Kokkos, RAJA, etc.<br><br>• Flang testing and evaluation<br><br>• Kitsune and Tapir | • Increasing dependence on LLVM<br><br>• Many vendors import and redistribute LLVM<br><br>• Contributions and collaborations with many vendors through LLVM<br>• AMD<br>• ARM<br>• Cray<br>• HPE<br>• IBM<br>• Intel<br>• NVIDIA |

*Other ECP activities with LLVM emerging organically.*

*Active involvement with broad LLVM community: LLVM Dev, EuroLLVM ECP personnel had 10+ presentations at the 2020 Dev Meeting*

# So how did we do?

# Reflections from ECP Panel at ECP AHM in May 2022

## Panel Overview

- As we enter the era of Exascale, it is time to reflect
  - What did we get right?
  - What did we miss?
  - What did we omit?
- Many of you participated in these workshops and reports
  - Comments and questions welcome!
  - Please use ZOOM Q&A window

Jeffrey Vetter (ORNL), Moderator
Pete Beckman (ANL)
Jack Dongarra (UTK, ORNL)
Bob Lucas (Ansys)
Kathy Yelick (UCB)

Vetter | ORNL

5

# 15 years is an eternity in computing - How did our predictions do?

## Hits

- System power came in at O(20MW) not O(1GW)
- Few major software rewrites / evolution
  - So far, FORTRAN -> C++ is the main conversion
- ECP included applications, software, and hardware
  - ~70 teams, ~1000 researchers
  - IESP
- Concurrency (1B-way parallelism)
- Open-source software

## Misses

- Systems deployed 4 years later than expected (of 2018)
- Programming systems are multiplying and immature/incomplete
- Hardware diversity
- Demise of vendor interest in HPC
- Fault tolerance



HOW STANDARDS PROLIFERATE:
(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC)

SITUATION: THERE ARE 14 COMPETING STANDARDS.

14?! RIDICULOUS! WE NEED TO DEVELOP ONE UNIVERSAL STANDARD THAT COVERS EVERYONE'S USE CASES. YEAH!

SOON:
SITUATION: THERE ARE 15 COMPETING STANDARDS.

## Overlooked

- Productive programming models (ala AI/ML): Python, Jupyter, Julia
- Cost of ECP + NRE + Procurements approaches $3.6B USD
- AI/ML is not predicted (or even mentioned)
- Cloud deployment models
- Green/sustainable computing

**OAK RIDGE**
National Laboratory

# Pondering Post-Exascale Computing

- Thinking about the next 10 years

# Important Architectural Trends

- Heterogeneous integration

- Chiplets

- Ecosystems and Standards
  - CXL, UCIe, BoW, …

- Open-source Tools and IP
  - RISC-V, OpenLane, Silicon Compiler, etc

- Open foundries

- *Codesign will be more important than ever*



CHIPS enables rapid integration of functional blocks at the chiplet level

Today — Monolithic    Tomorrow — Modular
Image: Intel

Custom chiplets    Commercial chiplets

COMM    RADAR EW    SIGINT

Figure 1. CHIPS Vision (DARPA)

[DARPA ERI Summit 2018]

AMD to Fuse FPGA AI Engines Onto EPYC Processors, Arrives in 2023
By Paul Alcorn published May 04, 2022

NVIDIA Opens NVLink for Custom Silicon Integration

New UCIe Chiplet Standard Supported by Intel, AMD, and Arm
By Paul Alcorn published March 02, 2022

Modular AMD Chips to Embrace Custom 3rd Party Chiplets
By Francisco Pires published June 17, 2022

University Shuttle Program
Spurring advanced semiconductor R&D

Intel Is Opening up Its Chip Factories to Academia
By Agam Shah

# Reimagining Codesign

- 2021 Workshop

- Four priority research directions

  - Drive Breakthrough Computing Capabilities with Targeted Heterogeneity and Rapid Design

  - Software and Applications that Embrace Radical Architecture Diversity

  - Engineered Security and Integrity from Transistors to Applications

  - Design with Data-Rich Processes

- We must make codesign agile, more accurate, and use real workloads



*Overview Brochure*

Basic Research Needs for

**Reimagining Codesign for Advanced Scientific Computing**

Unlocking Transformational Opportunities for Future Computing Systems for Science

16-18 March 2021

https://doi.org/10.2172/1822198

U.S. DEPARTMENT OF ENERGY | Office of Science

https://www.osti.gov/biblio/1822198-reimagining-codesign-advanced-scientific-computing-unlocking-transformational-opportunities-future-computing-systems-science

# Abisko: Microelectronics Codesign

# Abisko Microelectronics Codesign Overview



**Applications**

*Motifs, Composition*

**Algorithms**

*API, Motifs*

**Software**

*ISA, IR*

**Architecture**

*Circuit scale up, Interconnects, PDK*

**Devices and Circuits**

*Compact models*

**Materials**

Codesign

OAK RIDGE
National Laboratory

# Abisko Microelectronics Codesign Overview



1. Develop better techniques for codesign from algorithms to devices and materials
2. Design Spiking Neural Network chiplet that can be integrated with contemporary computer architectures
3. Explore new devices and materials for the SNN chiplet (neuron, synapse, plasticity, etc.)
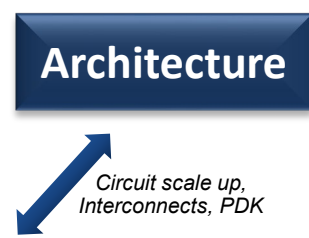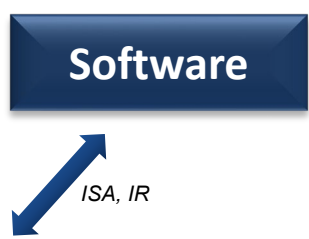4. Design language abstractions and runtime support for SNN chiplet



Source: *Wikipedia*

# Abisko Microelectronics Codesign Overview

OAK RIDGE National Laboratory

Sandia National Laboratories

ASU Arizona State University

Georgia Tech

HARVARD UNIVERSITY

Collaborator
Fermilab

1. Develop better techniques for codesign from algorithms to devices and materials
2. Design Spiking Neural Network chiplet that can be integrated with contemporary computer architectures
3. Explore new devices and materials for the SNN chiplet (neuron, synapse, plasticity, etc.)
4. Design language abstractions and runtime support for SNN chiplet

**Motivation**
- Transportation
- CMS Sensors

**Applications**

*Motifs, Composition*

CMS

F1Tenth



Source: *Wikipedia*

# Abisko Microelectronics Codesign Overview



Collaborator

**Fermilab**

1. Develop better techniques for codesign from algorithms to devices and materials
2. Design Spiking Neural Network chiplet that can be integrated with contemporary computer architectures
3. Explore new devices and materials for the SNN chiplet (neuron, synapse, plasticity, etc.)
4. Design language abstractions and runtime support for SNN chiplet

*Motivation*
- Transportation
- CMS Sensors

**Applications**

*Motifs, Composition*

**Algorithms**

*Algorithms*
- ML: SLAYER, Whetstone, EONS, eProp, STDP
- Non-ML: Graph algorithms, CSP
- Simulators: NEST, Brian2

*API, Motifs*



Source: *Wikipedia*

# Abisko Microelectronics Codesign Overview


(Oak Ridge National Laboratory, Sandia National Laboratories, ASU Arizona State University, Georgia Tech, Harvard University, Collaborator: Fermilab)

1. Develop better techniques for codesign from algorithms to devices and materials
2. Design Spiking Neural Network chiplet that can be integrated with contemporary computer architectures
3. Explore new devices and materials for the SNN chiplet (neuron, synapse, plasticity, etc.)
4. Design language abstractions and runtime support for SNN chiplet

**Motivation**
- Transportation
- CMS Sensors

*Motifs, Composition*

**Applications**

**Algorithms**
- ML: SLAYER, Whetstone, EONS, eProp, STDP
- Non-ML: Graph algorithms, CSP
- Simulators: NEST, Brian2

*API, Motifs*

**Algorithms**

**Software**
- DSL and API for neuromorphic co-processing
- Built on LLVM and MLIR
- Portable across Abisko chiplet, GPUs, etc.

*ISA, IR*

**Software**

*Source: Wikipedia*

# Abisko Microelectronics Codesign Overview

Oak Ridge National Laboratory

Sandia National Laboratories

ASU Arizona State University

GT Georgia Tech

HARVARD UNIVERSITY

Collaborator
Fermilab

1. Develop better techniques for codesign from algorithms to devices and materials
2. Design Spiking Neural Network chiplet that can be integrated with contemporary computer architectures
3. Explore new devices and materials for the SNN chiplet (neuron, synapse, plasticity, etc.)
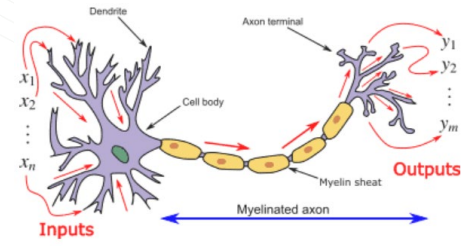4. Design language abstractions and runtime support for SNN chiplet

**Applications**

*Motivation*
- Transportation
- CMS Sensors

*Motifs, Composition*

CMS

*Algorithms*
- ML: SLAYER, Whetstone, EONS, eProp, STDP
- Non-ML: Graph algorithms, CSP
- Simulators: NEST, Brian2

**Algorithms**

nest::

BRIAN

F1Tenth

*API, Motifs*

*Software*
- DSL and API for neuromorphic co-processing
- Built on LLVM and MLIR
- Portable across Abisko chiplet, GPUs, etc.

**Software**

LLVM COMPILER INFRASTRUCTURE

MLIR

XACC

*ISA, IR*

*Architecture*
- Design neuromorphic chiplet
- RISC-V neuromorphic extensions
- Heterogeneous integration with contemporary technologies

Source: *Wikipedia*

**Architecture**

RISC-V

2.5D and 3D integration

Simulation/Emulation

ALADDIN

gem5

*Circuit scale up, Interconnects, PDK*

# Chiplet Architectures

- Design an (analog) SNN chiplet that can be easily integrated with contemporary technologies
  - Heterogeneous integration with mixed processes
  - Compatible with existing processes

- Extensive advances in chiplets, packaging, and heterogeneous integration recently
  - Open Domain-Specific Architecture
  - UCIe, BoW, TSMC SoIC-CoW, Intel Foveros

- Using open toolchain and architecture to explore chiplet designs: RISC-V, OpenLane



CHIPS enables rapid integration of functional blocks at the chiplet level

Today — Monolithic    Tomorrow — Modular

Image: Intel

Custom chiplets    Commercial chiplets

COMM    RADAR EW    SIGINT

| | Adaptive filter | | SerDes | | SerDes |
| | Beam forming | | Beam forming | | Adaptive filter |
| | QR Decomp. | | QR Decomp. | | QR Decomp. |

Figure 1. CHIPS Vision (DARPA)

[DARPA ERI Summit 2018]

Optional HBM DRAM Dies — HBM DRAM Die, TSV — Silicon interposer — Optional multiple logic dies

μBumps — Base Die — PHY — PHY — Compute — Logic

C4 Cu Bumps

Standard Package Trace

Package Balls

Package Substrate

Short Wires

Circuit Board

Figure 21. An example showing the use of 2D and 3D interconnections (courtesy TSMC)

[IEEE HIR 2021]

# Evaluation of 2.5D Chiplet for Neuromorphic Computing

Element Types:



Sandia VO$_2$ ECRAM

refine
improve

# ASIC Flow for Digital NN (baseline)

- Investigate the performance of fully customized ASIC design for ultra-fast NN inference
  - ORNL: HLS and RTL
  - Geogia Tech: ASIC Synthesis and PD
- Model details:
  - Fixed NN architecture with quantized weights
  - Experimented with 2bit or 3bit of inputs (limited by FermiLab implementation)
- Flow:
  - Vitis HLS to generate RTL
  - Catapult logic synthesis
  - Customized backend layout tool (incl. tech mapping, placement and routing)
- Achieved clock frequency of 1~ 2GHz in a 28nm technology



**conv 0**
**3bits   794 um**



**2bits   526 um**

# Abisko Microelectronics Codesign Overview



1. Develop better techniques for codesign from algorithms to devices and materials
2. Design Spiking Neural Network chiplet that can be integrated with contemporary computer architectures
3. Explore new devices and materials for the SNN chiplet (neuron, synapse, plasticity, etc.)
4. Design language abstractions and runtime support for SNN chiplet

*Motivation*
- Transportation
- CMS Sensors

**Applications**

*Motifs, Composition*

*Algorithms*
- ML: SLAYER, Whetstone, EONS, eProp, STDP
- Non-ML: Graph algorithms, CSP
- Simulators: NEST, Brian2

**Algorithms**

*API, Motifs*

*Software*
- DSL and API for neuromorphic co-processing
- Built on LLVM and MLIR
- Portable across Abisko chiplet, GPUs, etc.

**Software**

*ISA, IR*

*Architecture*
- Design neuromorphic chiplet
- RISC-V neuromorphic extensions
- Heterogeneous integration with contemporary technologies

**Architecture**

Source: *Wikipedia*

*Circuit scale up, Interconnects, PDK*

2.5D and 3D integration

Simulation/Emulation

MesaFAB ReRAM

*Devices and Circuits*
- ion insertion (reversible doping) sets analog states
- mRaman captures transition linear, non-linear switching
- Will extend to 36x36 x-bar array
- Electronic and other optical spectroscopies

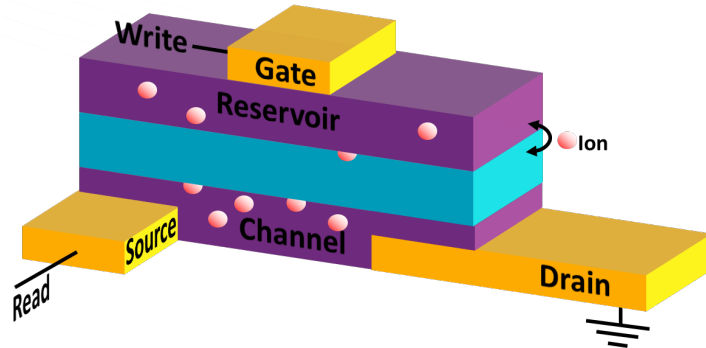**Devices and Circuits**

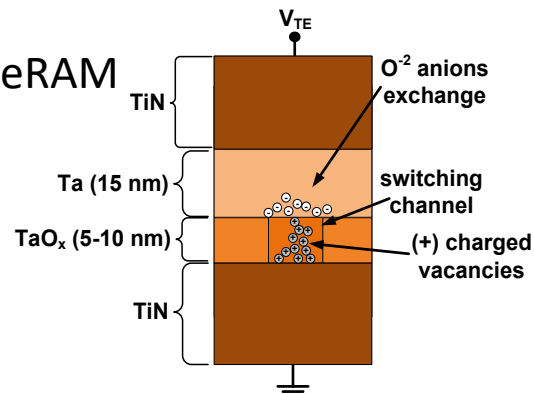*Compact models*

ECRAM

# Devices and Circuits

- Goals
  - Harness the interplay between mobile defects (ions and vacancies) and electronic properties to realize functional elements for spiking and non-spiking analog neuromorphic networks
  - Create and validate small network models; generate device and network data for co-design
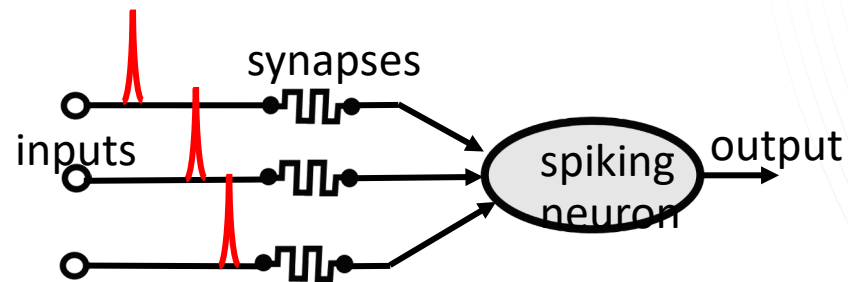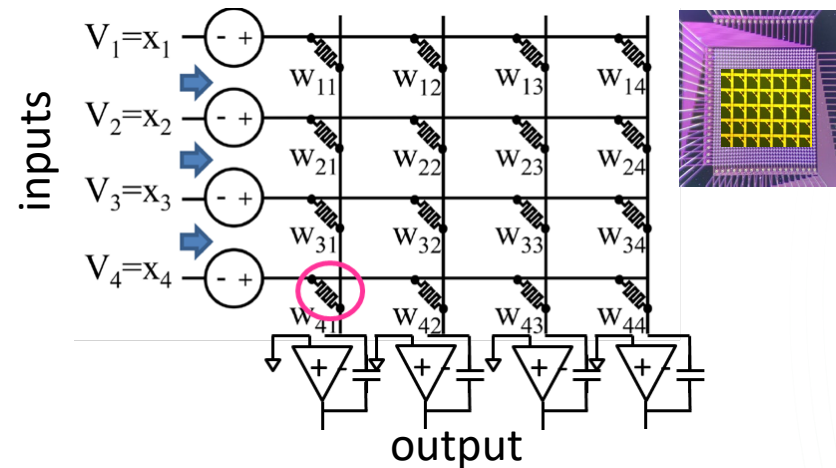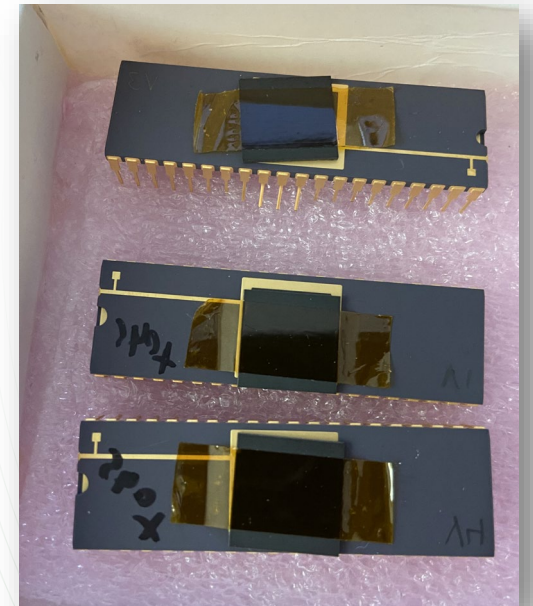  - Understand and mitigate radiation induced degradation mechanisms at the device and circuit level

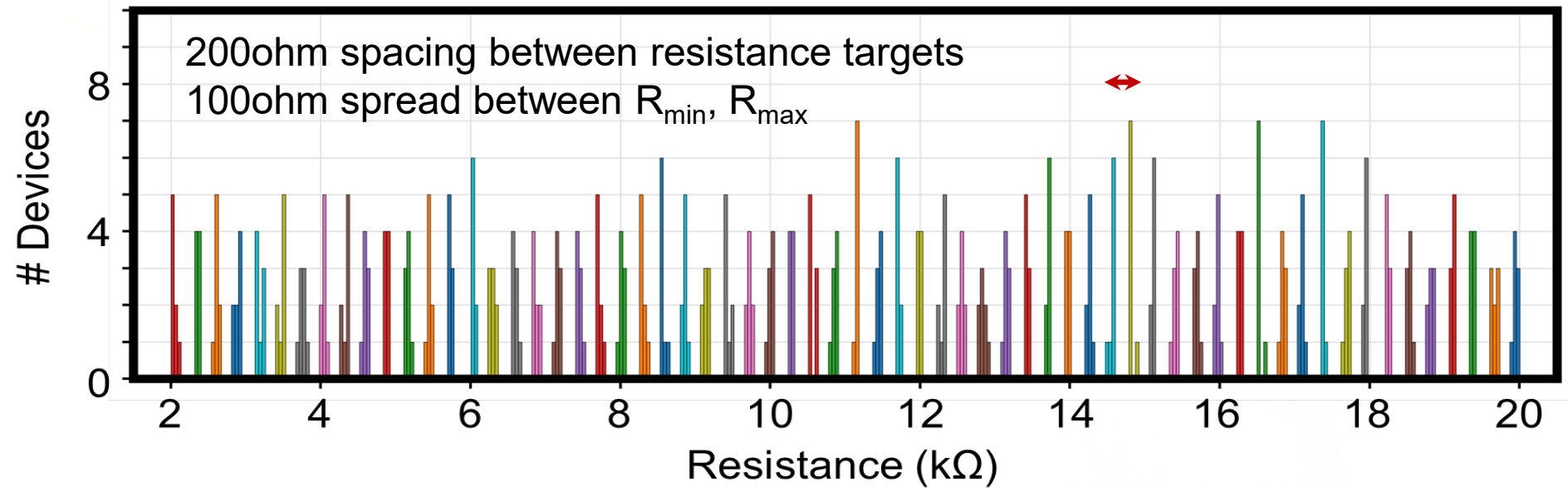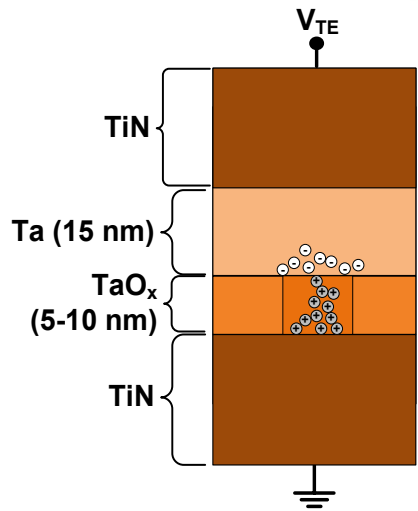# Experimental TaOx ReRAM Conductance Distributions

Developed TaOx weight mapping and programming routine for optimizing inference accuracy

# Abisko Microelectronics Codesign Overview

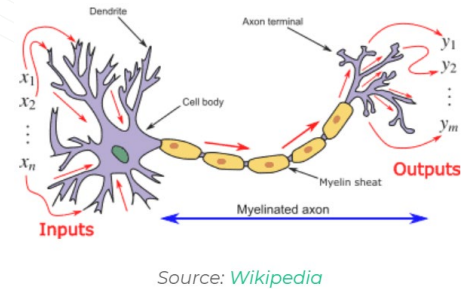OAK RIDGE National Laboratory
Sandia National Laboratories
ASU Arizona State University
GT Georgia Tech
HARVARD UNIVERSITY

Collaborator
Fermilab

1. Develop better techniques for codesign from algorithms to devices and materials
2. Design Spiking Neural Network chiplet that can be integrated with contemporary computer architectures
3. Explore new devices and materials for the SNN chiplet (neuron, synapse, plasticity, etc.)
4. Design language abstractions and runtime support for SNN chiplet

**Motivation**
- Transportation
- CMS Sensors

**Applications**

*Motifs, Composition*

CMS
3.8T Solenoid  ECAL  HCAL  IRON YOKE  Muon System Endcap (CSC+RPC)  TRACKER  21.6 m  15 m

**Algorithms**
- ML: SLAYER, Whetstone, EONS, eProp, STDP
- Non-ML: Graph algorithms, CSP
- Simulators: NEST, Brian2

**Algorithms**

*API, Motifs*

nest::
BRIAN

F1Tenth
VESC 6 MkIII Speed & Steering Control Board  Slamtec RPLIDAR A3  Power Distribution Board  UM7 IMU Board  Nvidia Jetson TX2  Traxxas Ford Fiesta Chassis & Drivetrain
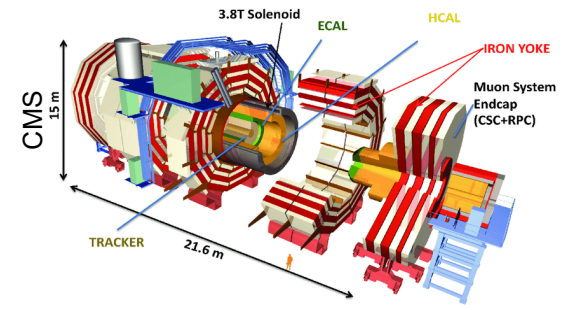
**Software**
- DSL and API for neuromorphic co-processing
- Built on LLVM and MLIR
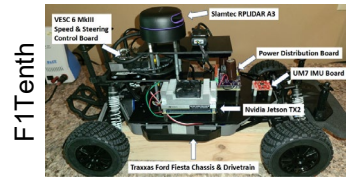- Portable across Abisko chiplet, GPUs, etc.

**Software**

LLVM COMPILER INFRASTRUCTURE
MLIR
XACC

*ISA, IR*

**Architecture**
- Design neuromorphic chiplet
- RISC-V neuromorphic extensions
- Heterogeneous integration with contemporary technologies

**Architecture**

RISC-V

Simulation/Emulation

2.5D and 3D integration
neuromorphic co-processor
memory for synaptic weights
Von Neumann main processor
inter-tier via
(a)

ALADDIN
gem5

Source: *Wikipedia*

Dendrite, Axon terminal, Cell body, Myelin sheath, Myelinated axon, Inputs, Outputs
$x_1$ $x_2$ $x_n$ $y_1$ $y_2$ $y_m$

**Devices and Circuits**
- ion insertion (reversible doping) sets analog states
- mRaman captures transition linear, non-linear switching
- Will extend to 36x36 x-bar array
- Electronic and other optical spectroscopies

**Devices and Circuits**

*Circuit scale up, Interconnects, PDK*

MesaFAB ReRAM

TaOx ReRAM
top metal (TiN), top ReRAM bit stack (tungsten), bottom via (tungsten), bottom metal (Al-Cu)

ECRAM
Write, Gate, Reservoir, Ion, Source, Channel, Drain, Read

phase 2, phase 1, electron probe, optical pump, network input, output
ROSS SIM

Computing Discovery Platform
neuron regime, synaptic regime, 40 mV/decade

*Compact models*

**Materials**
- Non-equilibrium probes to few nm
- Data-driven modeling
- On-demand neuromorphism

**Materials**

Domain wall memristor

Computational data mining

CNMS scanning probe microscopy and chemical imaging
bias-T, DC + 3 GHz, laser beam, AFM/PFM signal, sMIM-C, sMIM-G, $V_{bias}$, PZT, SRO, STO

Potentiation: $V_c$=-180 mV, $V_c$=-200 mV
Depression: $V_c$=-180 mV, $V_c$=-200 mV
ion insertion, ion removal

# Conclusions

# Recap

- Exascale is here!

- Our predictions were reasonably accurate, but we completely missed some
  - AI/ML
  - Programming systems remain major challenge

- Post-exascale
  - Heterogeneous integration and Chiplet architectures are vastly diversifying the architectural landscape
  - Post exascale will be accelerated by recent major semiconductor investments

- Abisko microelectronics codesign project developing a chiplet for analog SNN

- Start building your own chiplets today!

- Visit us (post COVID ☺)
  - We host interns and other visitors year round
    - Faculty, grad, undergrad, high school, industry

- Jobs at ORNL
  - Visit https://jobs.ornl.gov

- Contact me vetter@ornl.gov

- Experimental Computing Lab
  - Lots of emerging archs
  - https://excl.ornl.gov

OAK RIDGE
National Laboratory