# Abisko: Deep Codesign of an Energy-Optimized, High Performance Neuromorphic Accelerator

Jeffrey S. Vetter, ORNL (PI)

Alec Talin, Sandia NL

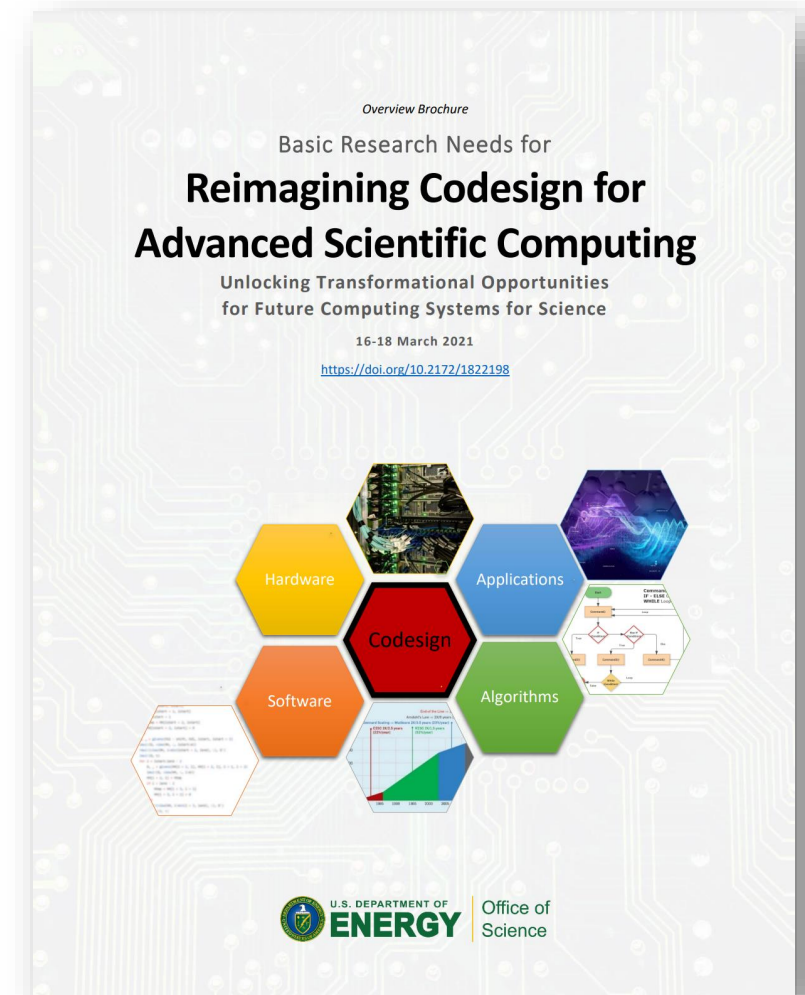David Brooks, Harvard

Yu Cao, ASU

Sung Kyu Lim, Georgia Tech

NICE Workshop
30 Mar 2022

**U.S. DEPARTMENT OF ENERGY**

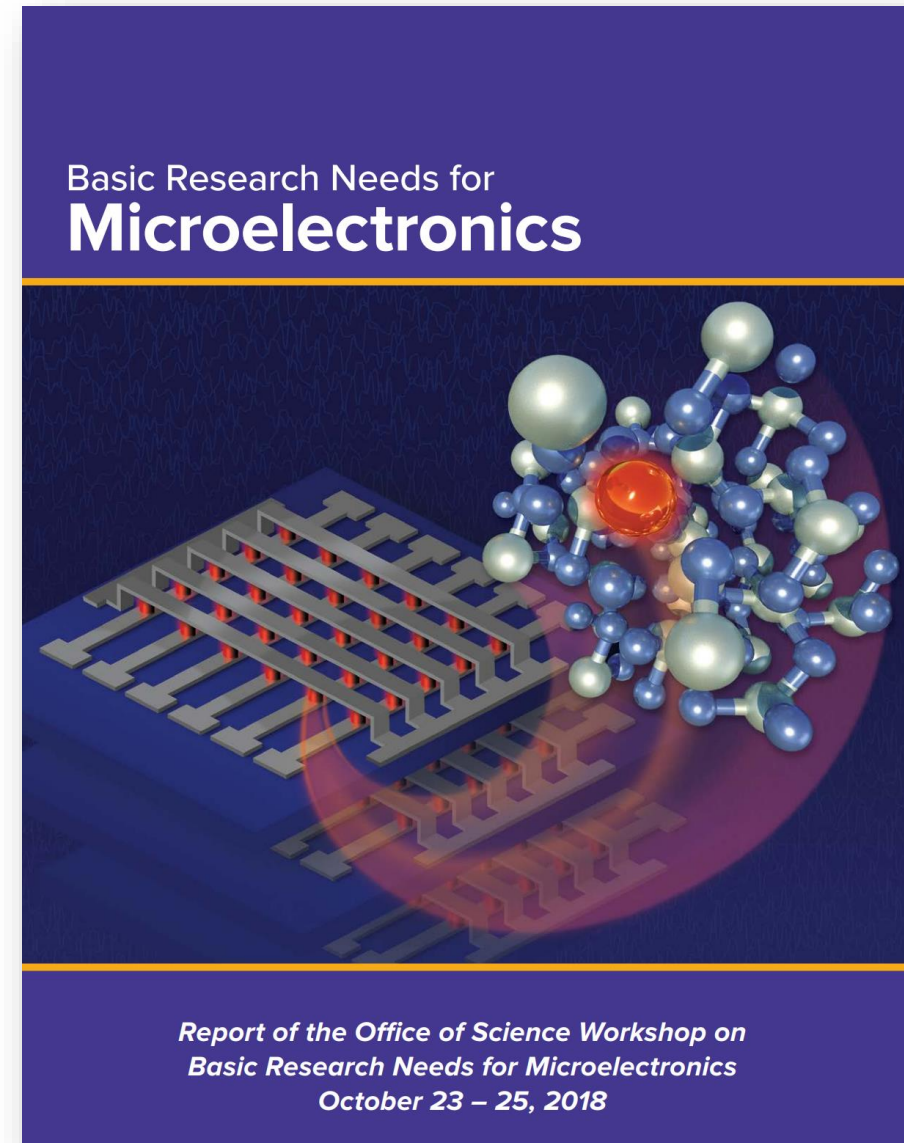ACSR Program Manager: Robinson Pino

# Overview

- Many factors are driving improved design of future computer systems
  - Electronics scaling, power, domain specific computing, business models, etc.
  - Massive demand for next-generation HPC systems (e.g., ModSim, AI, Data, Omniverse)

- DOE and others have embraced codesign as a path forward
  - Enable integrated design and implementation of end-to-end solutions, then iterate!
  - Reimagning Codesign focuses on new computational paradigms, workloads, agility

- Abisko is a new codesign project with the ambitious goals
  - Design Spiking Neural Network chiplet based on resistive switching materials that can be integrated with contemporary computer architectures
  - Develop portable software stack for neuromorphic algorithms across a range of platforms
  - Develop codesign framework for deep codesign into devices and materials

- Abisko is an interdisciplinary project including scientists from applications, algorithms, software, architectures, devices and circuits, and materials!



*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

# Basic Research Needs for Microelectronics (2018 Workshop)

- Five Priority Research Directions
  - <mark>Flip the current paradigm</mark>
  - Revolutionize memory and data storage
  - Reimagine informal flow unconstrained by interconnects
  - <mark>Redefine computing by leveraging unexploited physical phenomena</mark>
  - Reinvent the electricity grid through new materials, devices, and architectures



Basic Research Needs for **Microelectronics**

Report of the Office of Science Workshop on
Basic Research Needs for Microelectronics
October 23 – 25, 2018

# Recent DOE Program on Microelectronics Codesign



**Department of Energy**

## DOE Announces $54 Million for Microelectronics Research to Power Next-Generation Technologies

MARCH 24, 2021

Energy.gov » DOE Announces $54 Million for Microelectronics Research to Power Next-Generation Technologies

*National Labs Will Lead Transformation of Smart Devices, Clean Energy Technologies, and Semiconductor Manufacturing*

**WASHINGTON, D.C.** — The U.S. Department of Energy (DOE) today announced up to $54 million in new funding for the agency's National Laboratories to advance basic research in microelectronics. Microelectronics are a fundamental building block of modern devices such as laptops, smartphones, and home appliances, and hold the potential to power innovative solutions to challenges like the climate crisis and national security. Watch this video to learn more about microelectronics.
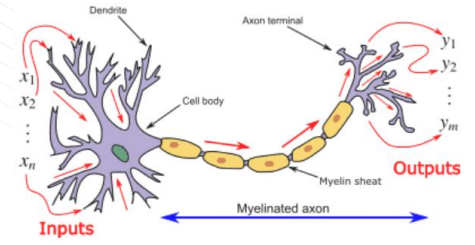
"Thanks to microelectronics, transformational technologies that used to swallow up entire buildings now fit in the palms of our hands—and it's time to take this work to the next level," said **Secretary of Energy Jennifer M. Granholm**. "Microelectronics are the key to the technologies of tomorrow, and with DOE's world-class scientists leading the charge, they can help bring our clean energy future to life and put America a step ahead of our economic competitors."

| Principal Investigator | Institution | City, State | Proposal Title |
|---|---|---|---|
| Guha, Supratik | Argonne National Laboratory (ANL) | Lemont, IL | Ultra-Dense, Near-Perfect, Atomic and Synaptic Memory |
| Taylor, Valerie | Argonne National Laboratory (ANL) | Lemont, IL | Threadwork: A Transformative Co-Design Approach to Materials and Computer Architecture Research |
| Braga, Davide | Fermi National Accelerator Laboratory (FNAL) | Batavia, IL | Hybrid Cryogenic Detector Architectures for Sensing and Edge Computing enabled by new Fabrication Processes |
| Garcia-Sciveres, Maurice | Lawrence Berkeley National Laboratory (LBNL) | Berkeley, CA | Co-Design and Integration of nano-sensors on CMOS |
| Ramesh, Ramamoorthy | Lawrence Berkeley National Laboratory (LBNL) | Berkeley, CA | Codesign of Ultra-Low-Voltage Beyond CMOS Microelectronics |
| Haegel, Nancy | National Renewable Energy Laboratory (NREL) | Golden, CO | Nitride materials and interfaces for radiation-hard integrated neutron detection |
| Vetter, Jeffrey | Oak Ridge National Laboratory (ORNL) | Oak Ridge, TN | Abisko: Codesign in the Wild: Designing Neuromorphic Hardware, Software, and Applications Concurrently using AI-enabled Methods |
| Graves, David | Princeton Plasma Physics Laboratory (PPPL) | Princeton, NJ | Diamond co-doping for quantum sensor applications |
| Aimone, James | Sandia National Laboratories (SNL) | Albuquerque, NM | COINFLIPS: CO-designed Improved Neural Foundations Leveraging Inherent Physics Stochasticity |
| McIntyre, Paul | SLAC National Accelerator Laboratory | Menlo Park, CA | Atoms-to-Systems Co-Design: Transforming Data Flow to Accelerate Scientific Discovery |

# Abisko Microelectronics Codesign Overview

**Collaborator**

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

1. Design Spiking Neural Network chiplet based on resistive switching materials that can be integrated with contemporary computer architectures
2. Develop portable software stack for neuromorphic algorithms
3. Extend codesign framework for deep codesign into devices and materials

*Source: Wikipedia*

**Applications**

*Motivation*
- Transportation
- CMS Sensors

*Motifs, Composition*

CMS

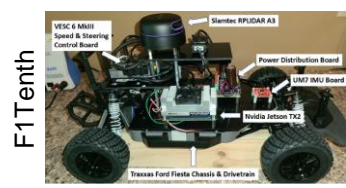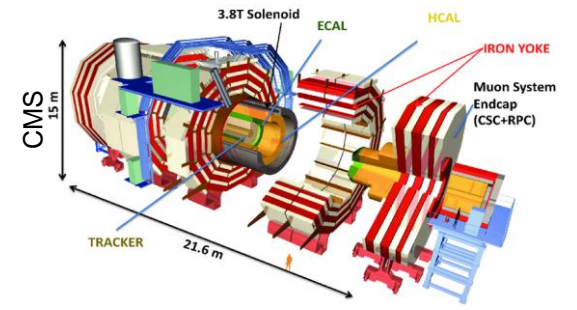*Algorithms*
- ML: SLAYER, Whetstone, EONS, eProp, STDP
- Non-ML: Graph algorithms, CSP
- Simulators: NEST, Brian2

**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

Simulation/Emulation
ALADDIN gem5

*ISA, IR*

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

**Architecture**

RISC-V

2.5D and 3D integration

MesaFAB ReRAM

*Circuit scale up, Interconnects, PDK*

*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**

TaOx ReRAM

ECRAM

ROSS SIM

Computing Discovery Platform

*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

# Team

Oak Ridge National Laboratory · Sandia National Laboratories · Arizona State University · Georgia Tech · Harvard University

- Vetter, Jeffrey S.
- Talin, Albert Alec (Devices)
- Kevin Cao
- David Brooks
- Lim, Sung-Kyu
- Comish, John
- Catherine "Katie" Schuman (Algo)
- Date, Prasanna
- Tripathy, D
- Farah Fahim
- Ghawaly, James
- Tallada, Marc Gonzonlas (Software)
- Gu-Yeon Wei
- Holland Hysmith
- Hornick II, Michael

- Huber, Joseph
- Ievlev, Anton
- Kulkarni, Shruti
- Frank Liu (Arch)
- Maksymovych, Petro (Materials)
- Marinella, Matthew
- Flynn, Michael
- Miniskar, Narasinga Rao
- Nhan Tran
- Ovchinnikova, Olga S.
- Sumpter, Bobby
- Aaron Young

# Abisko Microelectronics Codesign Overview

**Applications**

*Motivation*
- Transportation
- CMS Sensors

*Motifs, Composition*

**Algorithms**

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

*API, Motifs*

**Software**

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

*ISA, IR*

**Architecture**

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

*Circuit scale up, Interconnects, PDK*

**Devices and Circuits**

*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

*Compact models*

**Materials**

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

**CMS Experiment**
40MHz collision rate
~1B detector channels

Pb/s
40MHz

FPGA filter stack
~µs latency

10s Gb/s
~5 kHz

10s Tb/s
100s kHz

On-detector
ASIC compression
~100ns latency

On-prem CPU/GPU
filter farm
~100 ms latency

Worldwide
computing grid
Exabyte-scale
datasets

**1 Billion channels →**
**10x the average internet traffic in all of North**
**America**

# Pixel Detector: Proposed ML implementation

**Digital neuromorphic implementation**

Sensor (AFE) → Digital Neuron → Further compression

**Analog – Mixed Signal implementation using floating gates or memristive cross-bar arrays**

Sensor (Preamp) → On-chip feature classification and learning using reconfigurable network → ADC → Further compressive layers

- Ability to work in the latent space (downstream resources)

- Reconfigurability vs. pruning?

- On-chip inference vs. on-chip training?

- Light weight models?

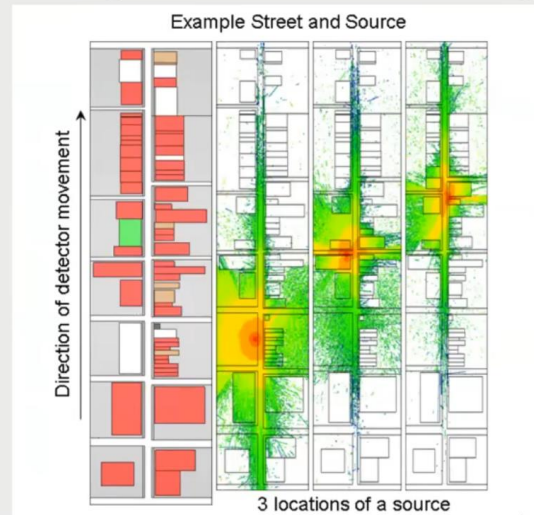- Can lead to self calibrating detectors?

# NeuroRad Project at ORNL

- 1: Develop a neuromorphic-capable radiation anomaly detection algorithm and evaluate on both simulated and real-world data.

- 2: Integrate neuromorphic algorithm on $\mu$Caspian board and integrate board with low power radiation detection system.
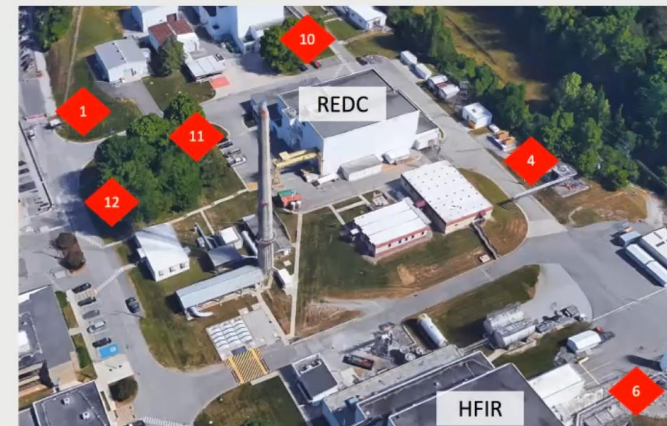
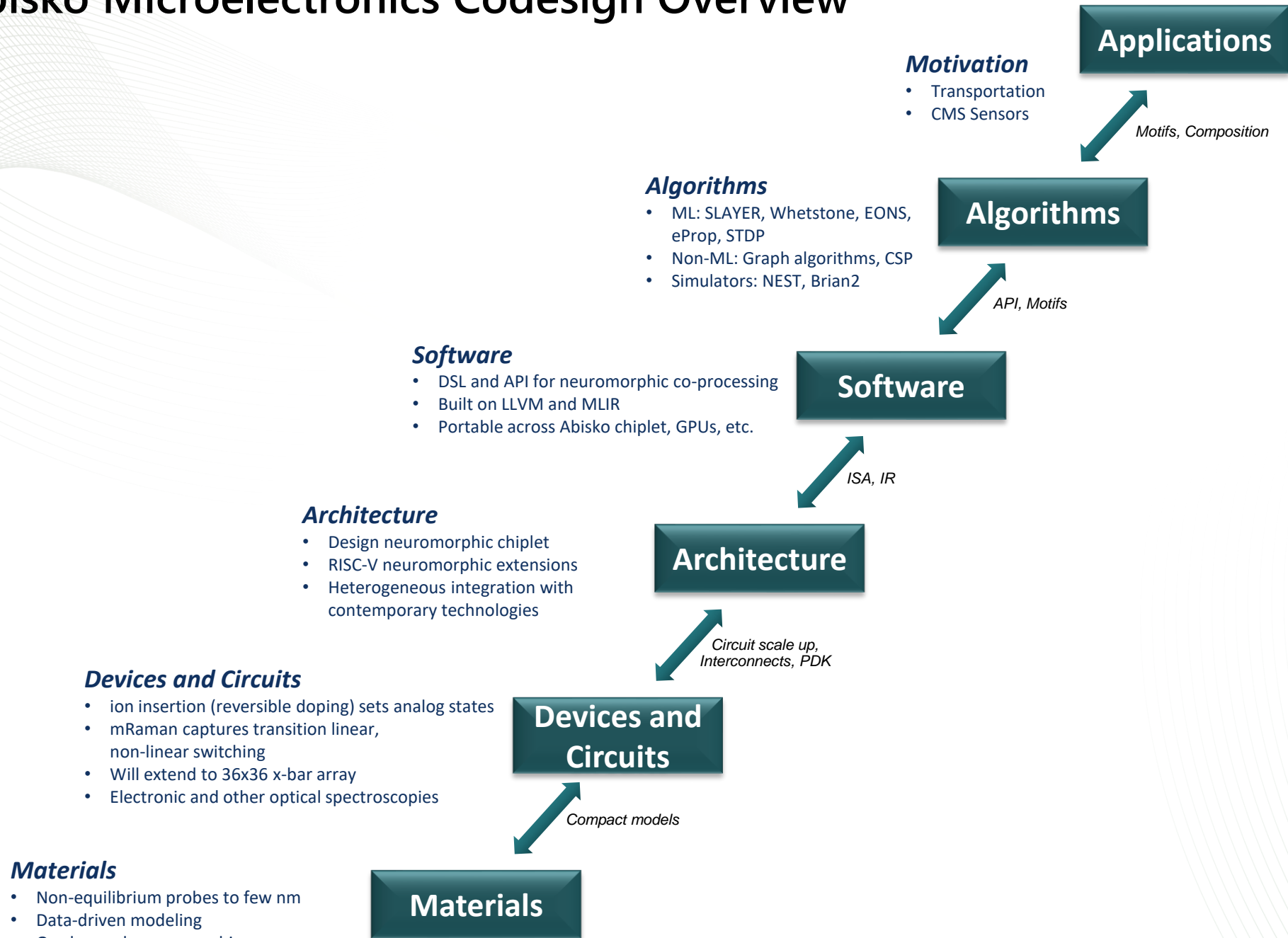| Datasets | |
|---|---|
| DOE Urban Search Challenge [1] | HFIR/REDC Static Monitors [2] |
| • Single 2"x4"x16" NaI(Tl) detector moving through urban street.<br>• 9700 training runs, 15840 testing runs | • Multiple static sensor "nodes" each with a single 2"x4"x16" NaI(Tl) detector, placed around ORNL HFIR/REDC facility.<br>• <200 source encounters |



Example Street and Source

Direction of detector movement

3 locations of a source



OAK RIDGE
National Laboratory

4

14

# Abisko Microelectronics Codesign Overview

**Applications**

*Motivation*
- Transportation
- CMS Sensors

*Motifs, Composition*

**Algorithms**

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

*API, Motifs*

**Software**

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

*ISA, IR*

**Architecture**

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

*Circuit scale up, Interconnects, PDK*

**Devices and Circuits**

*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

*Compact models*

**Materials**

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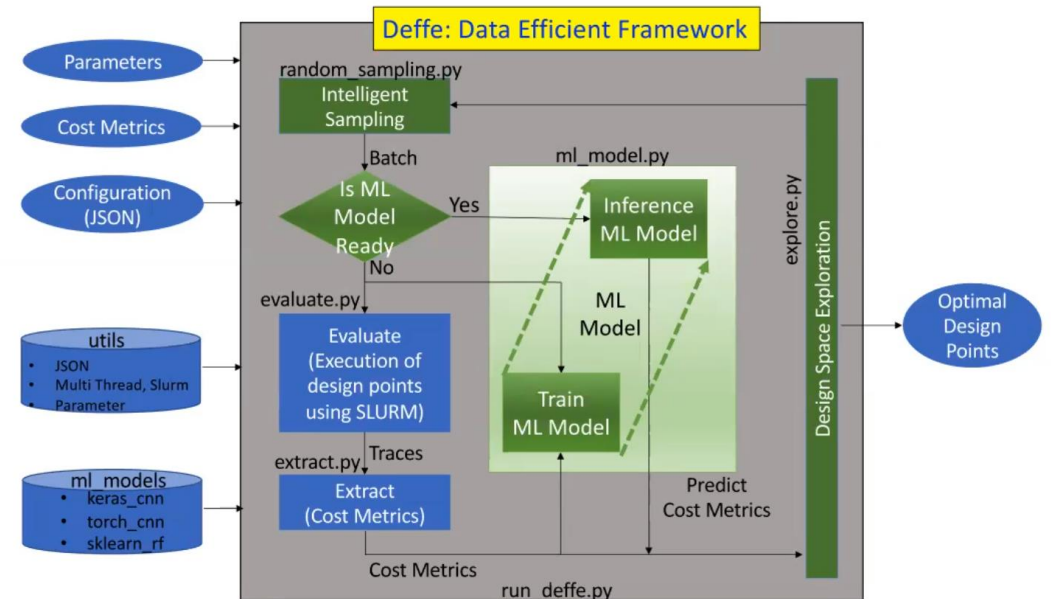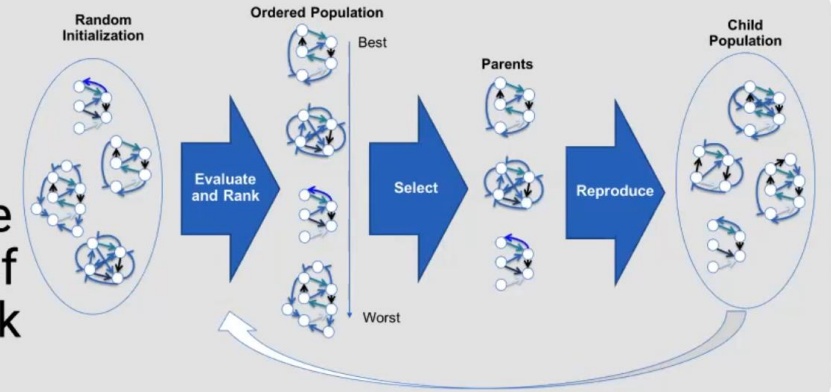- Data-driven modeling
- On-demand neuromorphism

# Algorithms

- Evaluate the best algorithms for specific problems
  - Include comparison against SOA techniques
- Evaluate algorithmic options for specific application
  - Input vector encoding
  - Evaluate different configurations with simulation
- Training, Inference, Online
- Interact with software and architecture teams
- Tools
  - EONS (Evolutionary optimization) for training
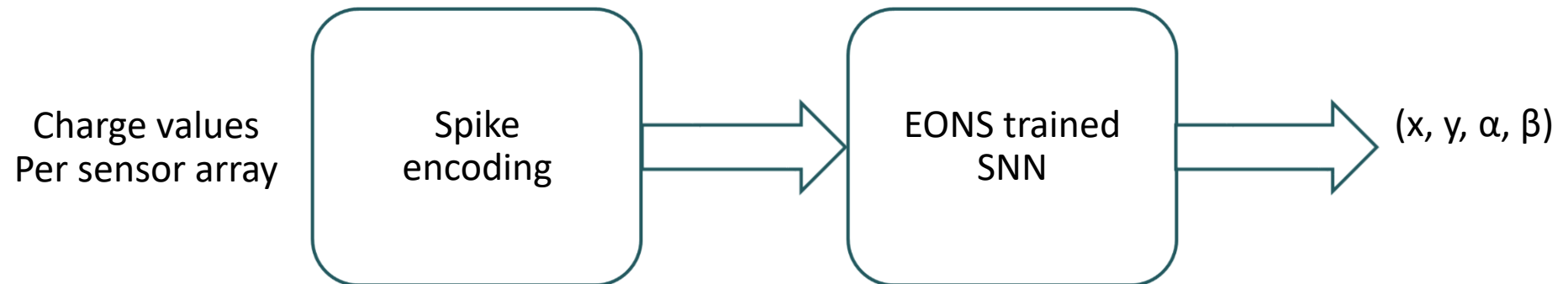  - Deffe for Hyperparameter optimization
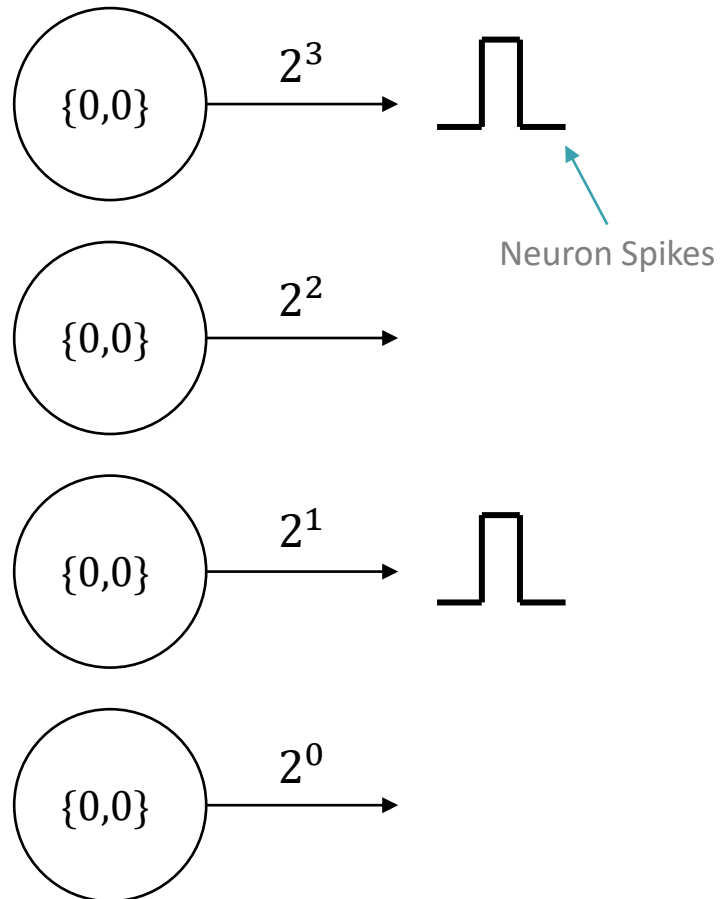
EONS

# Neuromorphic Approach for Smart Pixel Detection

- Dataset:
  - Charge values from the LHC every 250ps timesteps

- Goal
  - Data Compression, send only particle track information – (x, y, α, β)
  - In sensor pixel detection - hence, detection model needs to be small

- First approach
  - Apply neuromorphic algorithm – EONS
  - Explore spike encoding of charge values

- Other approaches
  - Regression, Spiking convolution NN, unsupervised learning (STDP), Spike-based Object detection algorithms
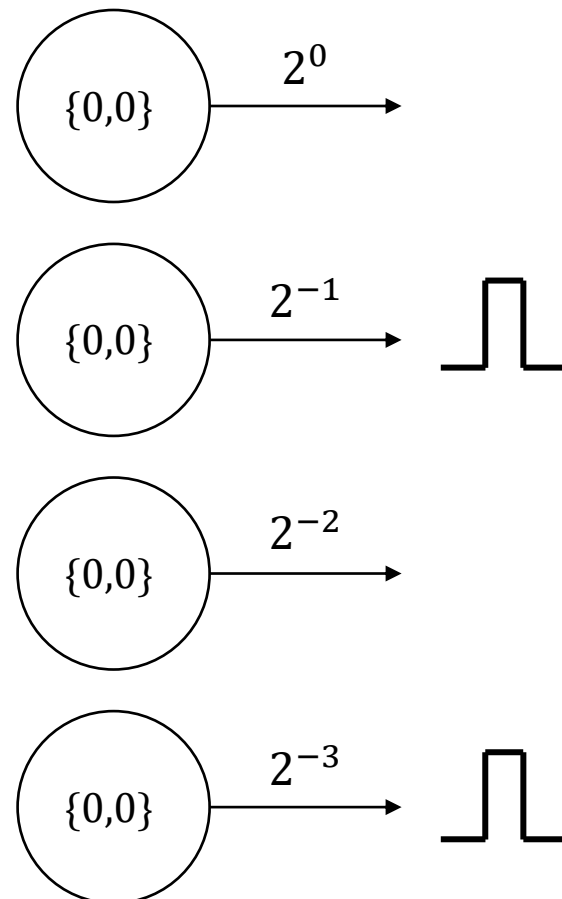
Charge values Per sensor array → [ Spike encoding ] → [ EONS trained SNN ] → (x, y, α, β)

# How to encode numbers on a neuromorphic computer?

Encoding 10, i.e. 1010

$2^3$

{0,0}

Neuron Spikes

$2^2$

{0,0}

$2^1$

{0,0}

$2^0$

{0,0}

Encoding 0.625, i.e. 0101

$2^0$

{0,0}

$2^{-1}$

{0,0}

$2^{-2}$

{0,0}

$2^{-3}$

{0,0}

Encoding -3.5, i.e. 0111

$-2^2$

{0,0}

$-2^1$

{0,0}

$-2^0$

{0,0}

$-2^{-1}$

{0,0}
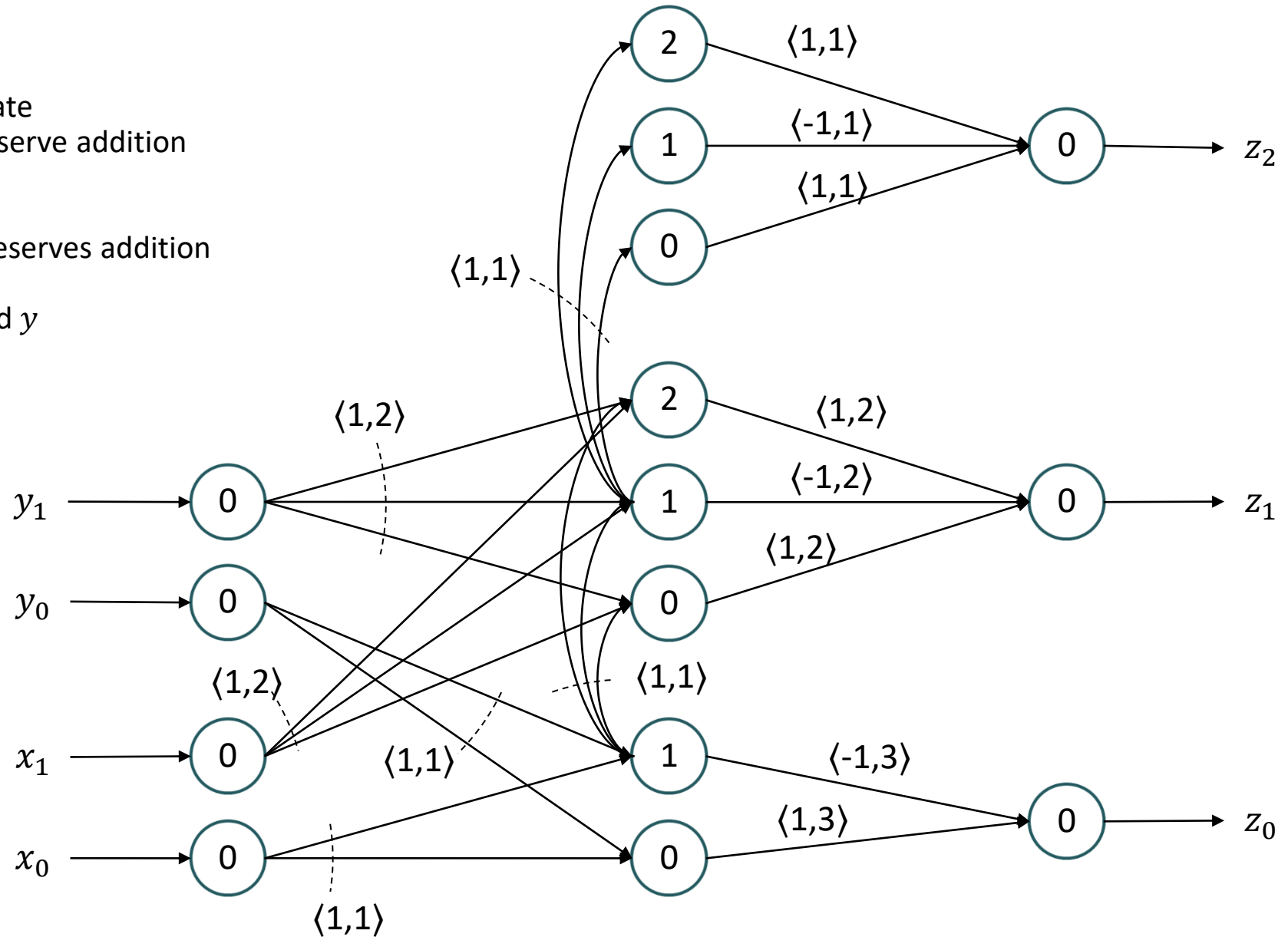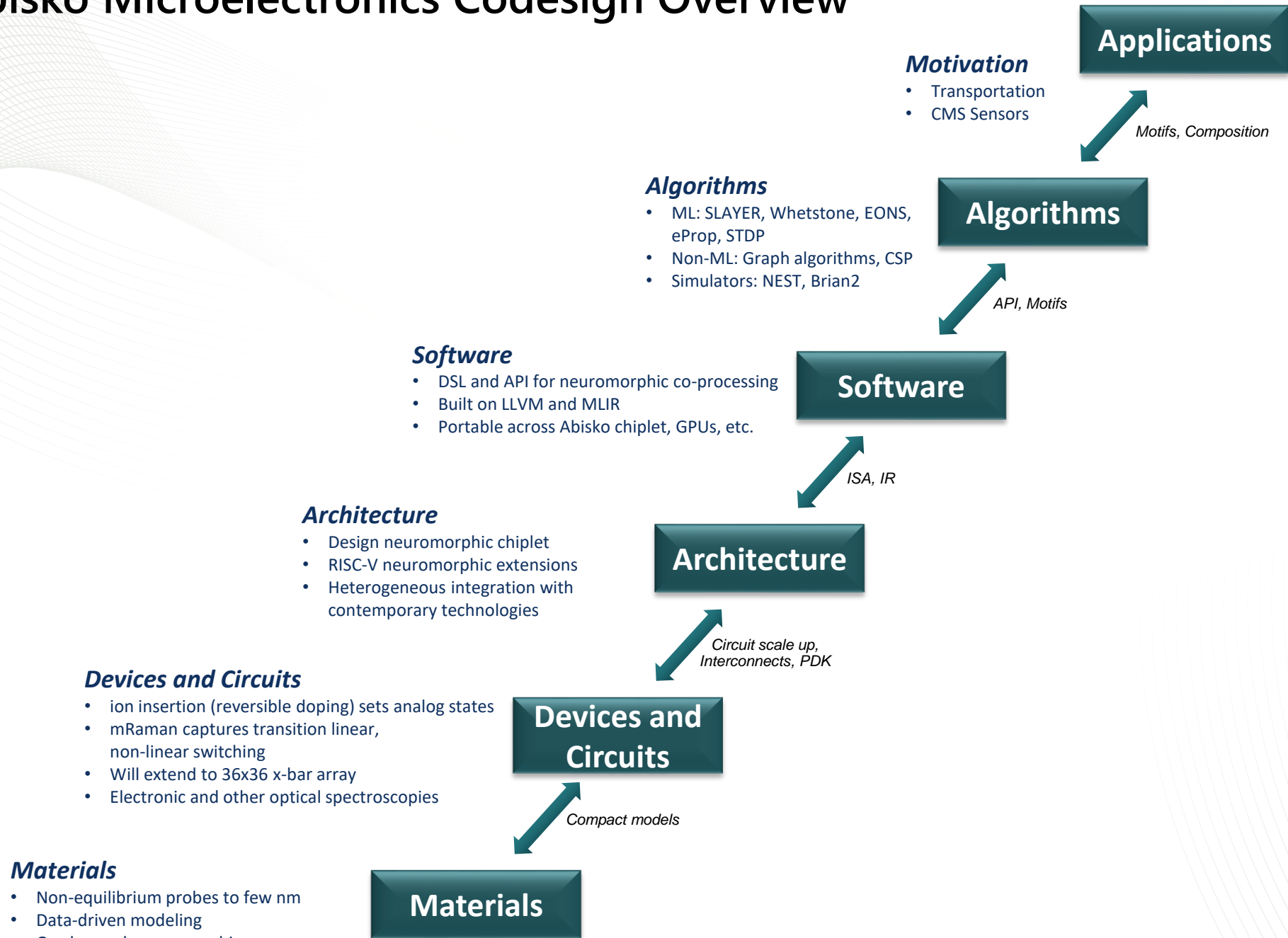
# The Virtual Neuron

- Current encoding methods are inadequate
  - Rate-based encoding does not preserve addition
  - Binning loses information

- Virtual neuron uses binary encoding, preserves addition

- Takes two 2-bit numbers as inputs: $x$ and $y$

- Returns a 3-bit number as output: $z$

- Implemented in NEST simulator

| $x_1$ | $x_0$ | $y_1$ | $y_0$ | $z_2$ | $z_1$ | $z_0$ | Sum |
|-------|-------|-------|-------|-------|-------|-------|-----|
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1+1=2 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 | 1+3=4 |
| 1 | 0 | 1 | 1 | 1 | 0 | 1 | 2+3=5 |
| 1 | 1 | 1 | 1 | 1 | 1 | 0 | 3+3=6 |

$y_1$ — 0

$y_0$ — 0

$x_1$ — 0

$x_0$ — 0

⟨1,1⟩  ⟨-1,1⟩  ⟨1,1⟩ → 0 → $z_2$

⟨1,2⟩  ⟨-1,2⟩  ⟨1,2⟩ → 0 → $z_1$

⟨-1,3⟩  ⟨1,3⟩ → 0 → $z_0$

⟨1,1⟩  ⟨1,2⟩  ⟨1,1⟩

# Abisko Microelectronics Codesign Overview

**Applications**

*Motivation*
- Transportation
- CMS Sensors

*Motifs, Composition*

**Algorithms**

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

*API, Motifs*

**Software**

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

*ISA, IR*

**Architecture**

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

*Circuit scale up, Interconnects, PDK*

**Devices and Circuits**

*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

*Compact models*

**Materials**

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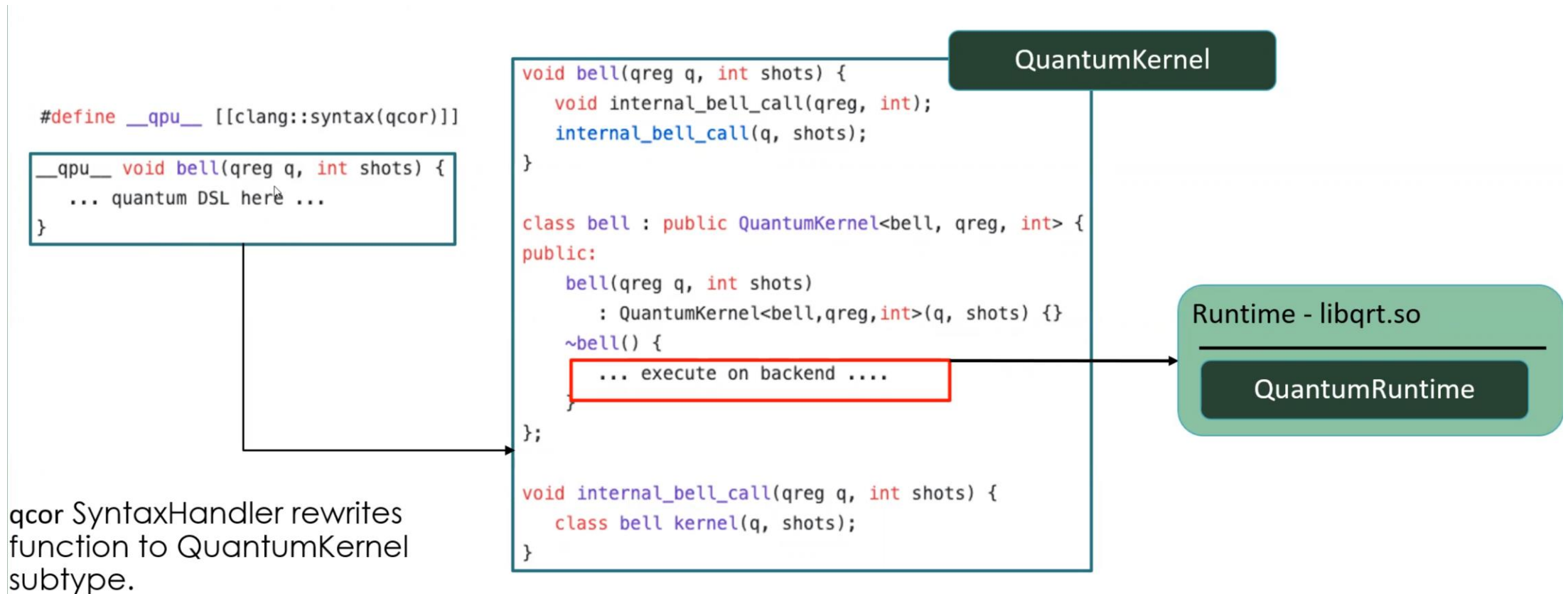- On-demand neuromorphism

# Software

- Develop a holistic software stack for neuromorphic coprocessing in a heterogeneous architecture
  - Programming model
  - Backend code generation
  - Runtime

- Portable to GPU, FPGA, SoC, and Abisko chiplet simulator

- Based on successful experiences with Quantum computing at ORNL:
  - XACC, QCOR

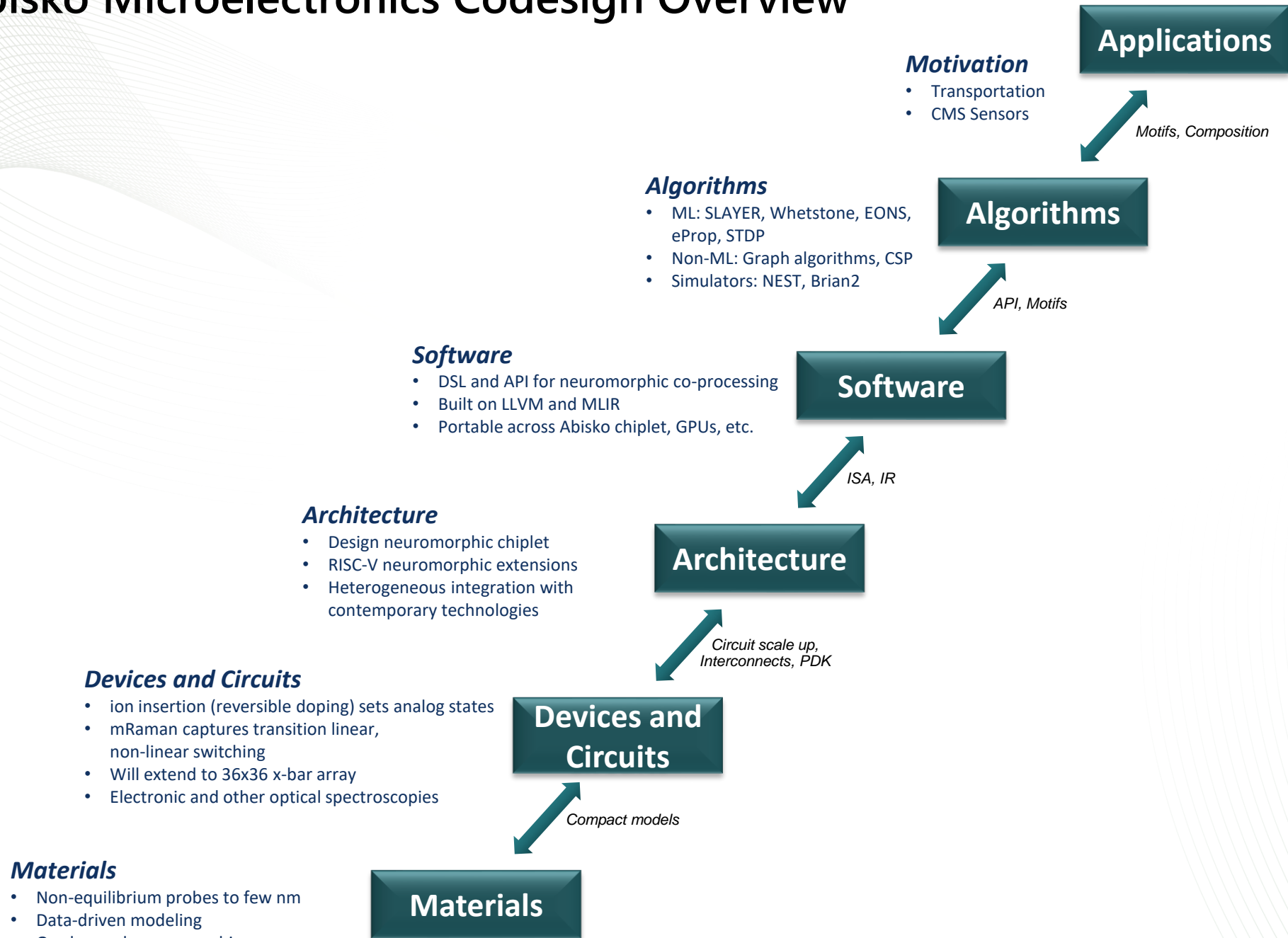- Building embedded DSL (Domain Specific Language) with LLVM and MLIR
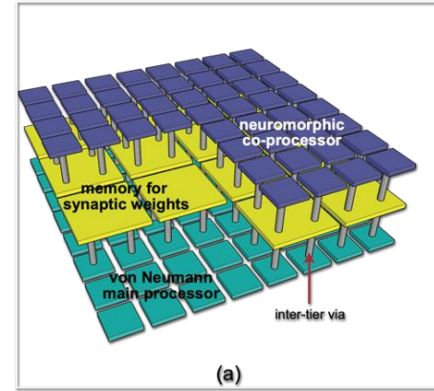
# XACC/QCOR Approach (as an analogue)



```
#define __qpu__ [[clang::syntax(qcor)]]

__qpu__ void bell(qreg q, int shots) {
    ... quantum DSL here ...
}
```

qcor SyntaxHandler rewrites function to QuantumKernel subtype.

```
void bell(qreg q, int shots) {
    void internal_bell_call(qreg, int);
    internal_bell_call(q, shots);
}


class bell : public QuantumKernel<bell, qreg, int> {
public:
    bell(qreg q, int shots)
        : QuantumKernel<bell,qreg,int>(q, shots) {}
    ~bell() {
        ... execute on backend ....
    }
};


void internal_bell_call(qreg q, int shots) {
    class bell kernel(q, shots);
}
```

QuantumKernel

Runtime - libqrt.so

QuantumRuntime

Program call to bell function is a call to another internal function that instantiates a temporary instance of the new QuantumKernel sub-type.

# Abisko Microelectronics Codesign Overview

**Applications**

**Motivation**
- Transportation
- CMS Sensors

*Motifs, Composition*

**Algorithms**

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

*API, Motifs*

**Software**

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

*ISA, IR*

**Architecture**

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

*Circuit scale up, Interconnects, PDK*

**Devices and Circuits**

**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

*Compact models*

**Materials**

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

# Architectures

- Design chiplet for SNN that can be easily integrated with contemporary technologies
    - Heterogeneous integration
    - Compatible with existing processes

- Use RISC-V interface to chiplet

- Simulate/emulate with existing simulators like Gem5 and Aladdin

# Abisko Architecture: Technology Landscape

- Advanced packaging is clearly one of the main technology drivers of semiconductor scaling in the near future



**From 2.5D to 3D and 3D+**

- 10-100X improvement / generation in data speed and bandwidth density

[Intel, TSMC, ISSCC 2021; IEEE HIR, 2021]

Y. Cao, ASU



**Roadmap of 3D Packaging**

- From 2010 to 2030: bandwidth density (Gbps/mm$^{-3}$) from <10 to $10^9$, energy efficiency (pJ/bit) from >1 to 0.01

[IEEE HIR, 2021; IMEC, 2021]

Y. Cao, ASU

- Underlying technology is the main uncertainty for neuromorphic accelerator

|  | SPIKING | NON-SPIKING |
|---|---|---|
| DIGITAL | CMOS-friendly (Loihi)latency and energy constraints | traditional GPU/FPGA/NN accelerators |
| ANALOG | interface to the rest of world, repeatability | Interface repeatability |

# Abisko Architecture: Smart Pixel Driver

- CMS Experiment from FemiLab: Farah Fahim
  - 40 MHz collision rate (25ns latency)
  - ~1B detector channels
- Active ongoing effort to design customized ASIC for data acquisition and compression
- Active ongoing effort to establish POR ML method on particle trajectory reconstruction



- On-going effort:
  - Establish baseline specs in computing intensity required using POR ML method
  - Explore techniques to better meet other constraints (quantization with fewer bits, spiking neuromorphic models)
  - Investigate and define Interface between ML accelerator cores and von-Neumann cores

# Abisko Microelectronics Codesign Overview

**Applications**

**Motivation**
- Transportation
- CMS Sensors

*Motifs, Composition*

**Algorithms**

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

*API, Motifs*

**Software**

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

*ISA, IR*

**Architecture**

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

*Circuit scale up, Interconnects, PDK*

**Devices and Circuits**

**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

*Compact models*

**Materials**

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

# Devices and Circuits Overview

*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



**Devices**

1) ECRAM

2) ReRAM

**Circuits**

# Experimental TaOx ReRAM Conductance Distributions

**Developed TaOx weight mapping and programming routine for optimizing inference accuracy**



200ohm spacing between resistance targets
100ohm spread between $R_{min}$, $R_{max}$
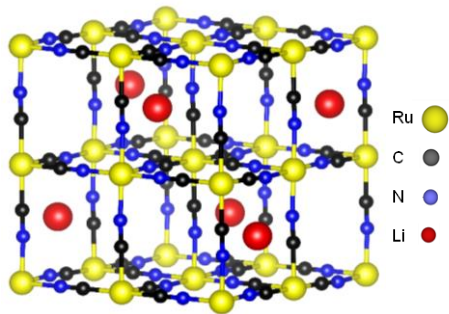
Resulting conductance distribution

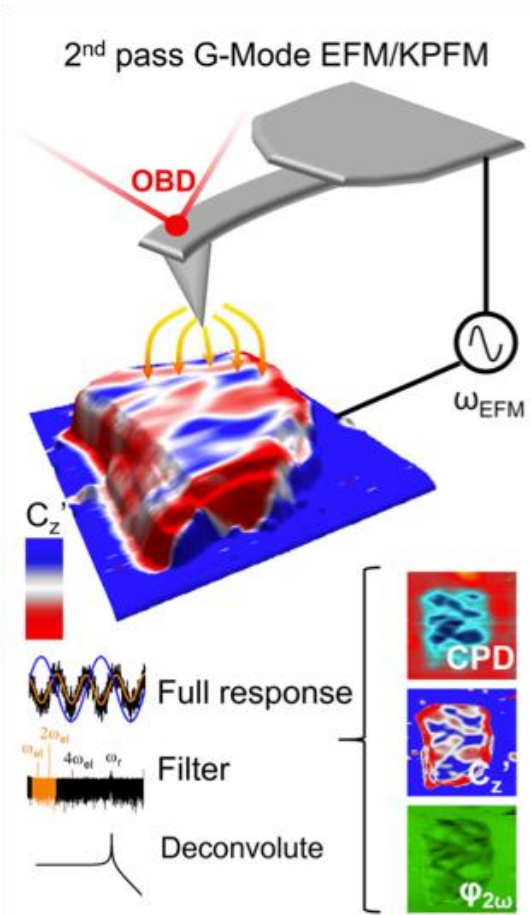# ECRAM Synapse Based on Ru Prussian Blue Analog

**Background**

- Prussian blue analogs are highly stable and can be patterned using photolithography
- Open crystal structure ideal for fast ion motion, but most PBA are poor electrical conductors

**Our work**

- We fabricate Ru PBA ECRAM synapses that switch with $Li^+$ or $H^+$ ions.
- The synapse display linear, symmetrical characteristics with excellent endurance and good retention
- Scaling experiments indicate $\Delta t_{sw} \sim 1ns$, $\Delta E_{sw} \sim 0.7fJ$ for 100 nm channel device.
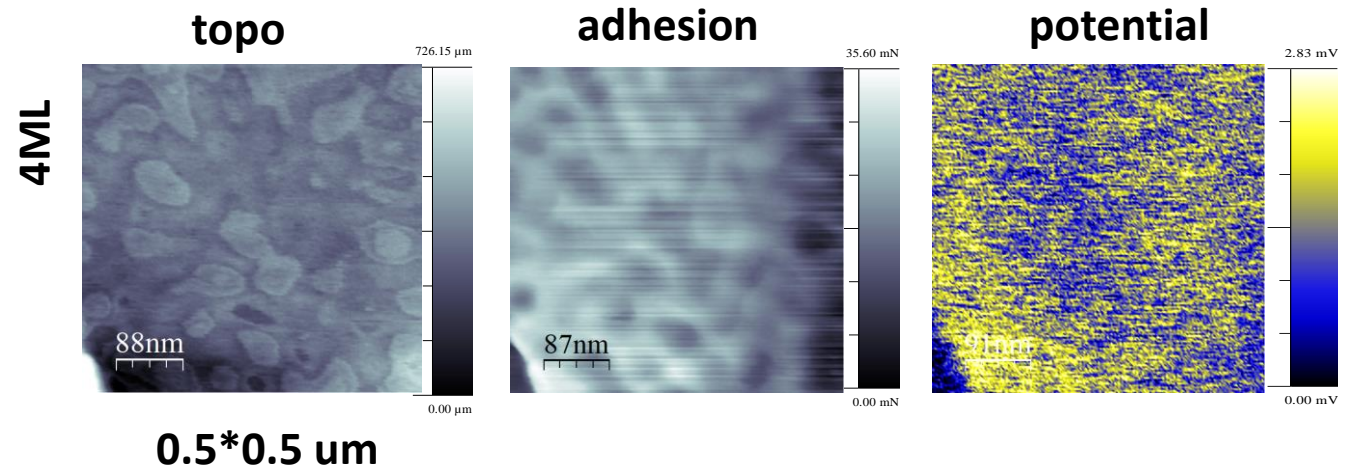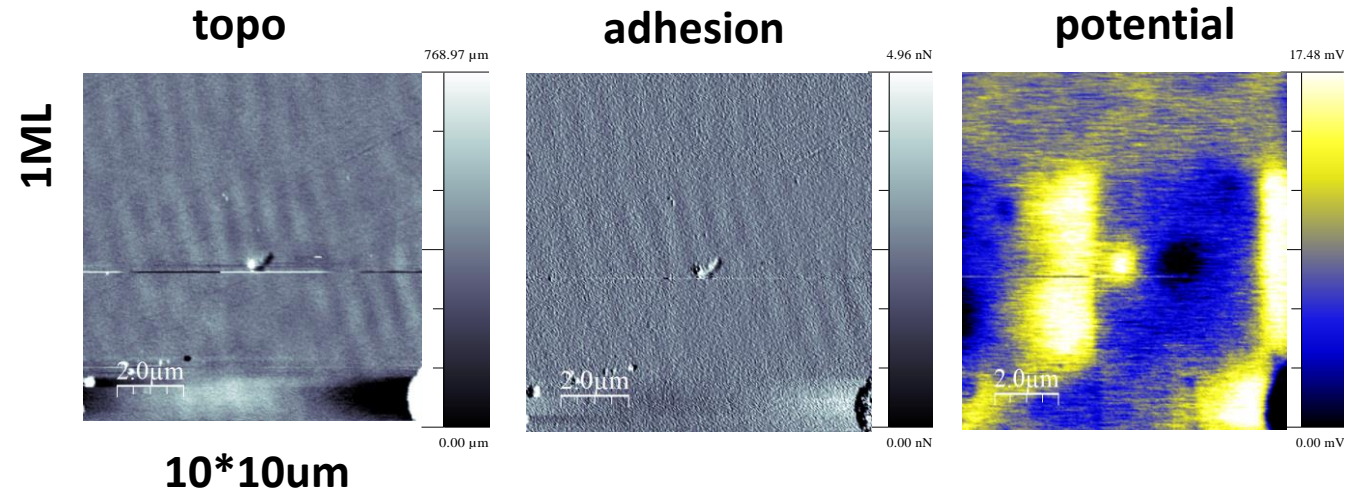
# Kelvin Probe Force Microscopy (KPFM) on PB thin films



The principles of the measurement procedure in KPFM technique using two pass mode

**M.Checa** et al, APL , **2021**



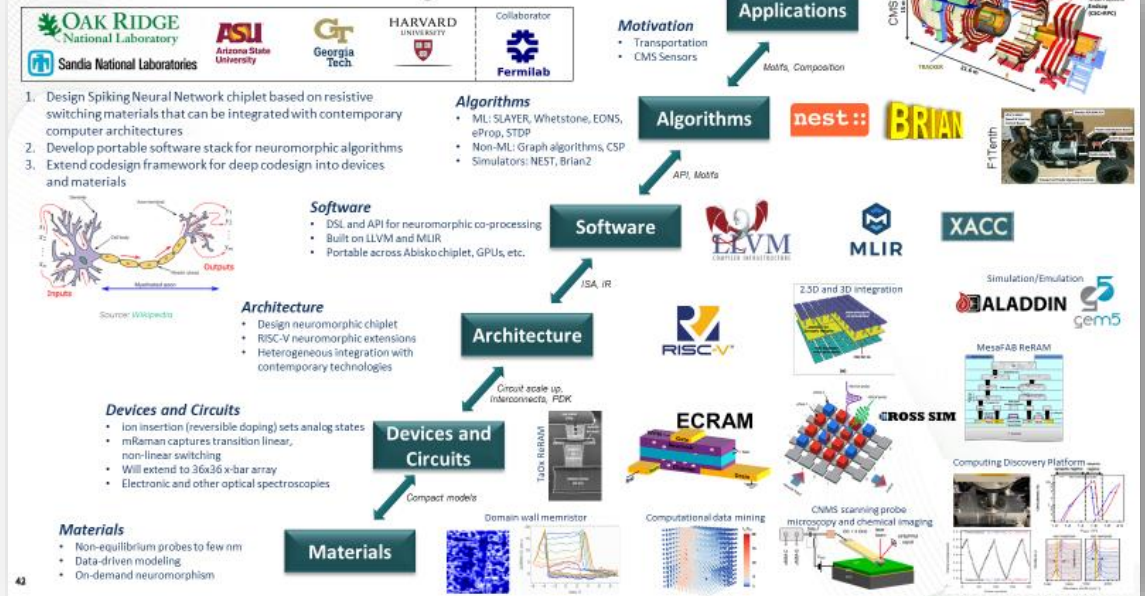Next step: nanoscale ionic effects from dielectric spectroscopy

# Devices and Circuits: Next steps and Recent Publications

- Investigate RuPBA fabrication compatible with Si integration
- Relate device characteristics to SPM measurements
- Develop compact model for ECRAM
- Construct small networks using TaOx memristors and Ru PBA elements
- Test radiation hardness, effects on accuracy, noise, and retention

- X. Xu, E. J.Cho, L. Bekker, **A. A. Talin**, E. Lee, A. J. Pascall, M. A. Worsley, J. Zhou, C. C. Cook, J. D. Kuntz, S. Cho and C. A. Orme, *A Bioinspired Artificial Injury Response System Based on a Robust Polymer Memristor to Mimic a Sense of Pain, Sign of Injury and Healing*, Adv. Science 2200629, 2022

- Su-in Yi, **A. A. Talin, M. J. Marinella**, R. S. Williams, *Physical Compact Model for Three-Terminal SONOS Synaptic Circuit Element*, submitted

# Summary



Abisko Microelectronics Codesign Overview

- Abisko is a new microelectronics codesign project with goals of
  - Design Spiking Neural Network chiplet based on resistive switching materials that can be integrated with contemporary computer architectures
  - Develop portable software stack for neuromorphic algorithms across a range of platforms
  - Develop codesign framework for deep codesign into devices and materials

- Truly interdisciplinary team working across the stack

- More information
  - vetter@computer.org
  - https://vetter.github.io

- We are hiring!
  - See jobs.ornl.gov
  - Send an email to me.

# Thanks!

**OAK RIDGE** National Laboratory

# Bonus Material