

# Advanced Computing, Sensing, and Algorithms for Highly Automated Driving

Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia LLC, a wholly owned subsidiary of Honeywell International Inc. for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.

## Meeting purpose

---

- Share and provide feedback on technology gaps and research plans for computing in highly-automated vehicles
- Identify useful metrics for energy-efficient computing and sensing
- Discussion topics
  - Sensing technologies
  - Low power, edge computing
  - Artificial intelligence and machine learning
  - Simulation and data

# National imperative



EXECUTIVE OFFICE OF THE PRESIDENT  
WASHINGTON, D.C.



July 31, 2018

M-18-22

MEMORANDUM FOR THE HEADS OF EXECUTIVE DEPARTMENTS AND AGENCIES

FROM: MICK MULVANEY  
DIRECTOR, OFFICE OF MANAGEMENT AND BUDGET

MICHAEL KRATSIOS  
DEPUTY ASSISTANT TO THE PRESIDENT  
OFFICE OF SCIENCE AND TECHNOLOGY POLICY

SUBJECT: FY 2020 Administration Research and Development Budget Priorities

“Agencies should prioritize investment in research and infrastructure to maintain U.S. leadership in strategic computing, **from edge devices to high-performance computing, that accelerates delivery of low power, high performance devices**; supports a national high-performance computing ecosystem; and explores novel pathways to advance computing in a post-Moore's Law era”.



“Today, semiconductors underpin the most exciting ‘must-win’ technologies of the future, including artificial intelligence to power self-driving cars and other autonomous systems...”

To secure America’s leadership in these future technologies for the next 50 years, the United States must continue to lead the world in semiconductor research, design, and manufacturing”

# National Strategic Computing Initiative



## NATIONAL STRATEGIC COMPUTING INITIATIVE UPDATE: PIONEERING THE FUTURE OF COMPUTING

*A Report by the*

FAST-TRACK ACTION COMMITTEE ON STRATEGIC COMPUTING

NETWORKING & INFORMATION TECHNOLOGY  
RESEARCH & DEVELOPMENT SUBCOMMITTEE

COMMITTEE ON SCIENCE & TECHNOLOGY ENTERPRISE

*of the*

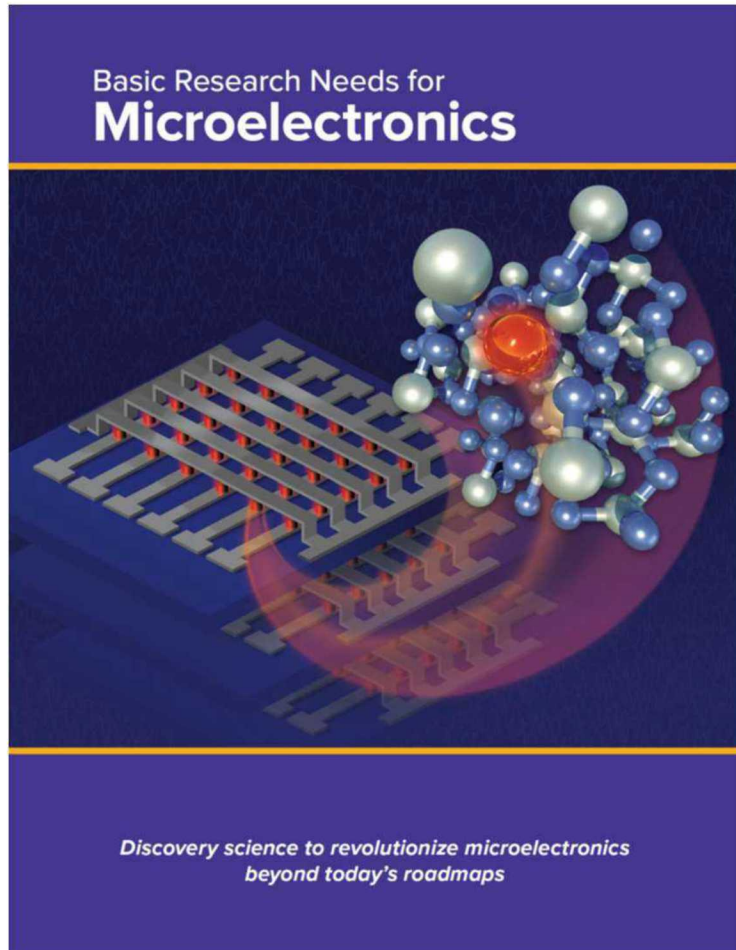
NATIONAL SCIENCE & TECHNOLOGY COUNCIL

NOVEMBER 2019

1. Pioneer new frontiers of digital and non-digital computation to address the scientific and technological challenges and opportunities of the 21st century.
2. Develop, broaden, and advance the Nation's computational infrastructure and ecosystem.
3. Forge and expand partnerships for the future of computing to ensure American leadership in science, technology, and innovation.



# DOE exploring low-energy electronics and advanced computing



Priority research directions:

1. Flip the current paradigm: define innovative material, device, and architecture requirements driven by applications, algorithms, and software
2. Revolutionize memory and data storage
3. Reimagine information flow unconstrained by interconnects
4. Redefine computing by leveraging unexploited physical phenomena

# DOE exploring low-energy electronics and advanced computing



**WORKSHOP ON ADVANCED COMPUTING FOR CONNECTED & AUTOMATED VEHICLES**

Date: May 7, 2019

The U.S. Department of Energy's (DOE) Vehicle Technologies Office (VTO) invites you to a Workshop on Advanced Computing for Connected & Automated Vehicles (CAVs) at Lawrence Berkeley National Laboratory in Berkeley, California.

This one-day summit will explore advanced microelectronics and computing approaches that can help meet future energy, cost, and computational requirements for CAVs. The workshop will bring together experts from the microelectronics industry, autonomous vehicle innovators, national laboratories, and academia in a precompetitive forum to discuss critical questions, including:

- What system sensing and computing architectures will fully automated vehicles require, and how much energy will those technologies consume?
- Which advanced computing approaches could reduce the energy requirements for fully automated vehicles while meeting their computational requirements?

**RSVP TODAY TO JOIN THE DISCUSSION**

<http://www.cvent.com/d/16q0h3>

**MAY 2019**

**7**

Lawrence Berkeley National Laboratory  
Berkeley, CA

U.S. DEPARTMENT OF **ENERGY**

In Cooperation With:

SLAC National Accelerator Laboratory  
DOE Office of Energy Efficiency & Renewable Energy  
DOE Vehicle Technologies Office  
DRAPER

Will highly automated vehicles be viable with conventional computing approaches, or will they require a step-change in computing?

What are the energy requirements to support on-board sensing and computing for highly automated vehicles?

What advanced computing approaches could reduce the energy requirements for highly automated vehicles while meeting their computational requirements?



# Projected computing performance and power



computing must  
meet size, weight,  
and power  
constraints

~1 petaflops  
~100 W (system)  
~100 TOPS/watt (SoC)



Full level 5 automated driving

TOPS == Trillion (tera) Operations



Early prototype self-driving

<https://www.wired.com/story/self-driving-cars-power-consumption-nvidia-chip/>

~100 teraflops  
~1000 W (system)  
~1 TOPS/watt (SoC)

>10x less  
power

>10x  
compute

>100x  
power  
performance

Significant innovation will be required in  
microelectronic materials and devices,  
sensing and computing architectures, and  
computer algorithms.

# Why now for the computing industry?

---

## Technology:

- End of Dennard power scaling; power becomes the key constraint
- Slow-down in Moore's Law, evidenced by flattening of transistor cost takedown

## Architectural:

- Limitation and inefficiencies in exploiting instructional-level parallelism and the prevailing von-Neumann architecture

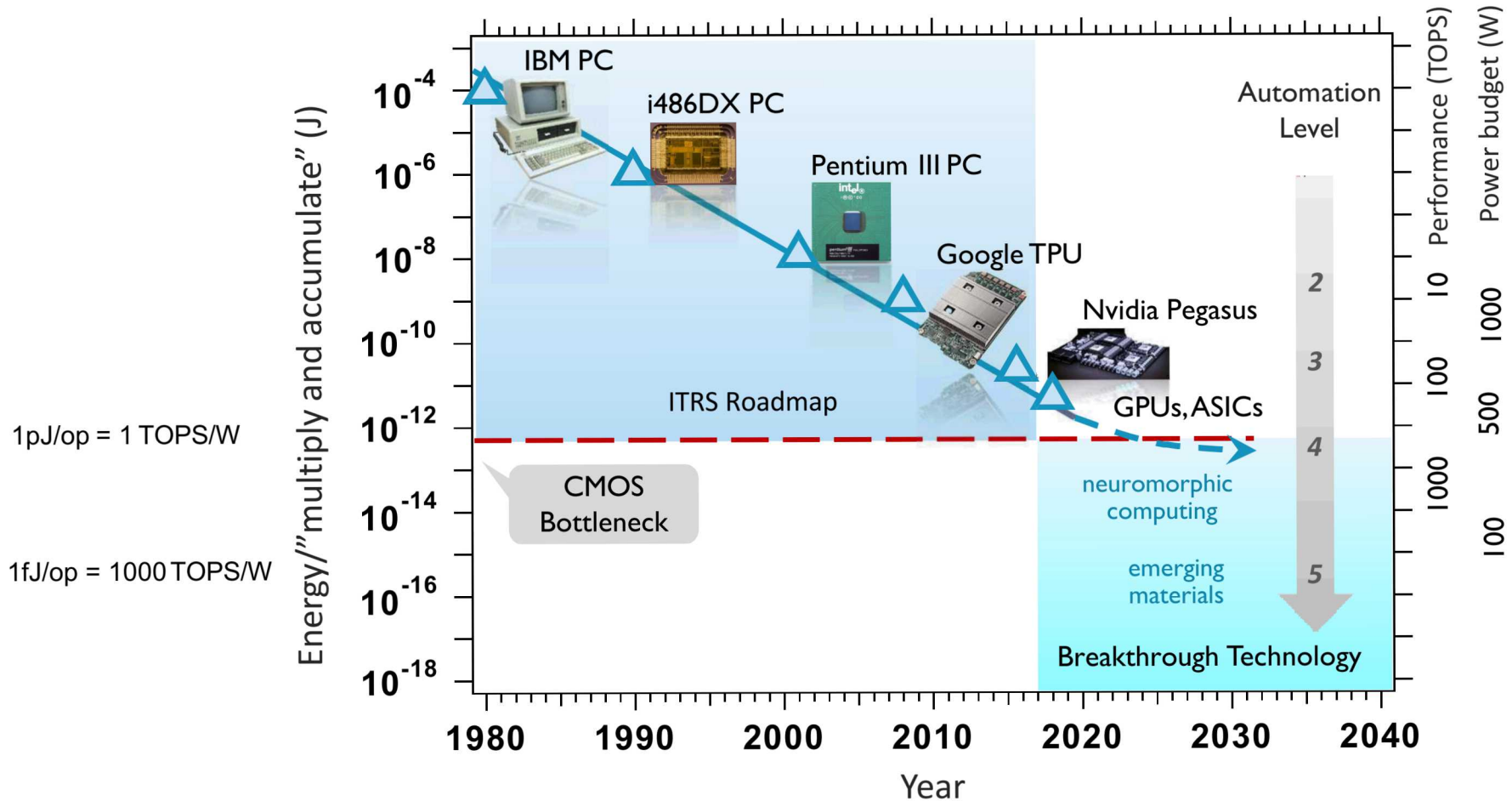
## Applications:

- Shift from desktop to mobile and IoT
- Ultra-scale cloud computing and artificial intelligence/machine learning workloads

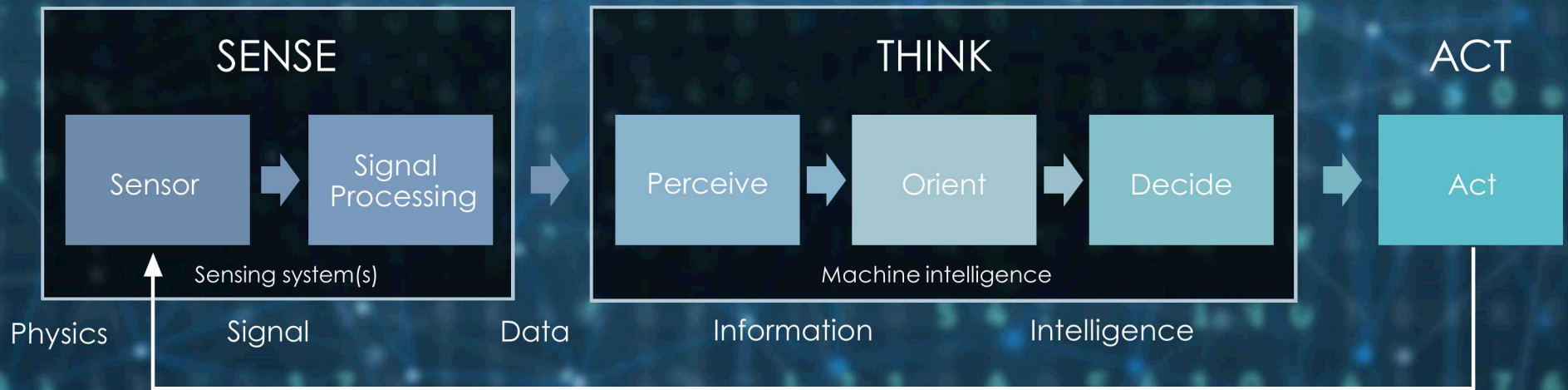
## Industry collaborations:

- End of International Technology Roadmap for Semiconductors (ITRS) roadmap
- Decline in SRC participation and the end of SEMATECH (absorbed in SUNY)

# Performance targets require breakthrough technology



# Automated systems





# Sandia's focused research areas for highly automated vehicles

Sensors and Sensor Processing	Scene Perception and Algorithms	Navigation Hardware Accelerators
<ul style="list-style-type: none"><li>• Create disruptive optical sensing technology to reduce energy consumption by 100X</li><li>• Develop chip-scale LiDAR to reduce cost by 100X</li></ul>	<ul style="list-style-type: none"><li>• Explore sparse coding and reduced-precision to reduce computation load by 1000X</li><li>• Develop biologically inspired machine learning algorithms to reduce the number of training samples by 100X</li><li>• Develop unsupervised and self-supervised learning algorithms</li></ul>	<ul style="list-style-type: none"><li>• Develop and demonstrate hardware capable of real-time processing of tera- to petabit inputs, with energies at &lt;10 fJ per operation (&gt;100 TOPS/W)</li><li>• Develop algorithms for robust and reliable recognition tasks needed for perception.</li><li>• Demonstrate the value of algorithm and hardware co-design such that combined elements have greater energy and/or SWaP improvement</li></ul>

## Leverage broad Sandia capabilities



COMBUSTION RESEARCH FACILITY



MESA MICROFAB



COMPUTING & INFORMATION  
SCIENCE



CENTER FOR INTEGRATED  
NANO TECHNOLOGIES



# Sandia Cooler



- Sandia Cooler technology has advanced through a DOE Technology Commercialization Fund project with industry partner Wakefield-Vette; now at TRL 8 with partner Heico
- Technology demonstrated in solid-state lighting for commercial warehouse applications
  - LED are located on rotating frame, ~1000 W power inductively coupled
  - Approximately 500 W of heat rejection
- Idea for CAV computing cooling – embed computing devices on rotating frame (similar to lighting) and communicate with adjacent vehicle data streams through 5G wireless link

## Meeting purpose

---

- Share and provide feedback on technology gaps and research plans for computing in highly-automated vehicles
- Identify useful metrics for energy-efficient computing and sensing
- Discussion topics
  - Sensing technologies
  - Artificial intelligence and machine learning
  - Low power, edge computing

# Sensor modalities for highly automated driving

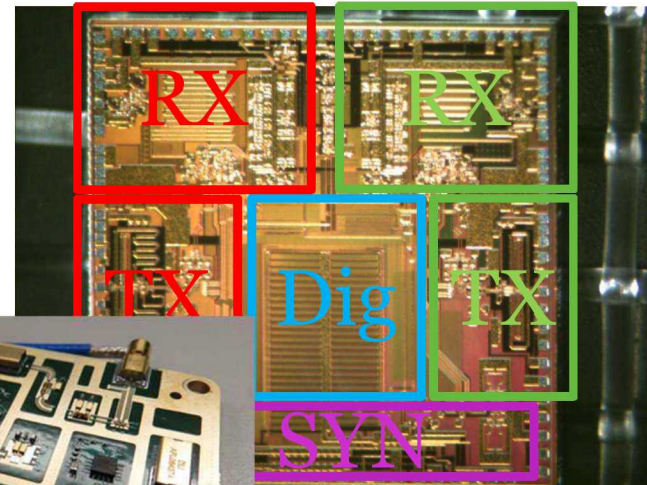
# Sensor design, integration, and data interpretation expertise

- Sandia has decades of navigational expertise
  - Radio Frequency/Acoustic (GPS, radar, sonar)
  - Inertial Navigation (accelerometers)
  - Guidance Systems (telemetry, tracking algorithms)
- Designed, fabricated, and deployed navigation components
  - Radar systems/RF Microwave Components
  - Gyroscopes (laser ring) and 6 axis accelerometers
  - Imagers (X-ray, optical, radar)
- Pioneered radar image processing and precision GPS-denied navigation
  - Unique Algorithm development
  - High consequence computationally intense image processing and real-time object recognition and tracking

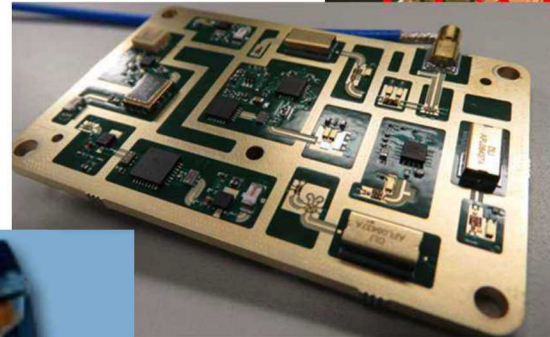
Mini Synthetic Aperture Radar



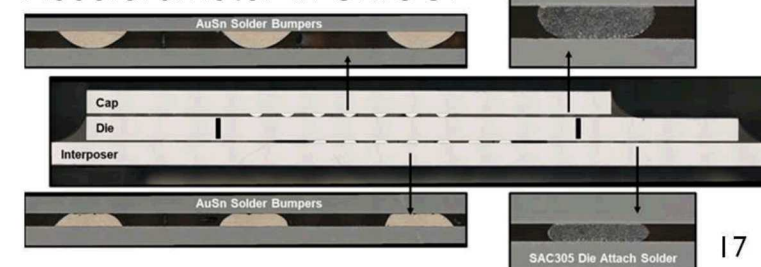
X and Ku Dual Band Radar RFIC



RX TX Module



Accelerometer in CMOS7





# Imaging radar

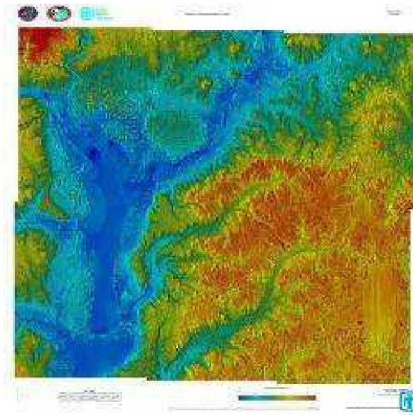
- Excellent complement to other sensors (electro-optical and LiDAR) for automated driving
  - Operates through all weather (fog, rain, snow)
  - Self illuminating (day, night)
  - Electrical scanning doesn't require moving parts
- High resolution, optical-like
  - Existing bands (76-81 GHz) provide centimeter class resolution
  - High frequency sensors result in small antennas/components
  - Resolution is not dependent on range to target
- Favorable computation complexity
  - Moving object detection (position/velocity vector) is a native product of radar, low computational complexity
  - Full radar image formation is computationally expensive but not needed in automotive applications
  - Image processing has significantly less computational cost than other imaging modalities
  - **Important because sensor data processing for useful information dominates complexity!**



SpotDwell image of a building at Jacksonville Naval Air Station.

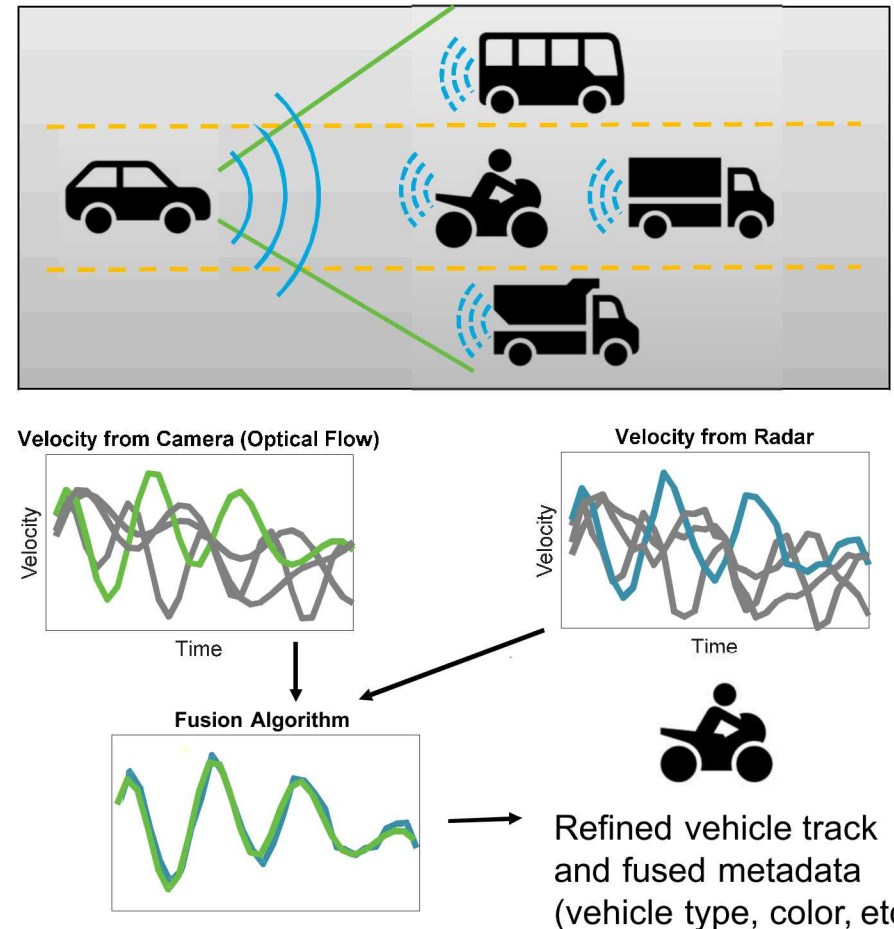
# SNL imaging radar heritage

- 35 years of experience building real-time, high-resolution, low SWaP radar imaging platforms
- Pioneered many new image processing and exploitation techniques – and continue to innovate with new algorithms and methods
- Experts in harsh environment high performance electrical systems
- Expertise in low power mm Wave RFIC component design



# Sensor data fusion

- Sensor fusion is the combination and exploitation of raw data at the sensor level rather than the derived data level
- Requires tight sensor integration
- Recent advancements at Sandia have been made in multi-sensor processing, but few multi-sensor platforms exploit true sensor fusion
- An example in the automotive arena would be the association of motion from moving vehicles in a **radar return** with object detected in **video**<sup>1</sup>
- The information resulting from sensor fusion will be higher confidence than the sum of information from sensors in isolation
- Allows application specific computing, in parallel to decision computing
- Sandia is a leader in adopting sub-threshold ASIC design (100x reduction in power)



[1] Naething, Richard M., and Richard C. Ormesher. "Doppler-assisted sensor fusion." U.S. Patent No. 10,267,895. 23 Apr. 2019.

# Chip-scale beam scanners



# LiDAR for the masses

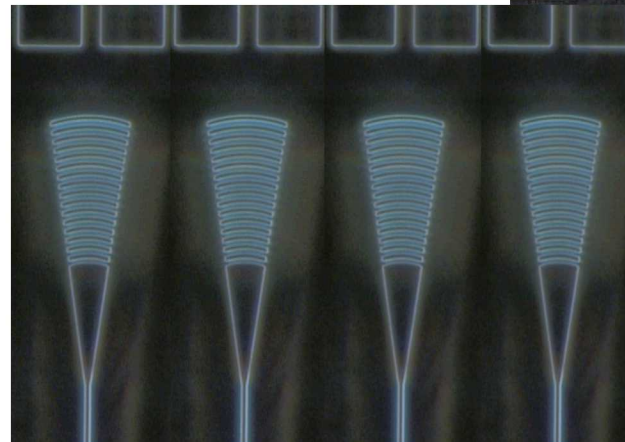
- Motivation:
  - autonomous vehicles require sensors that are robust, precise and easily manufacturable
  - current sensors use mechanical mirrors to steer light → large and expensive to produce
- Solution
  - Phased Arrays - recast established techniques at a new wavelength
  - Interfering waves create a narrow beam of light
  - Apply small electrical signals to adjust optical phase and steer beam
  - Requires precise fabrication at light wave dimensions and immense scalability



Very Large Array radio telescope –  
New Mexico



Mobile lidar mapping units  
atop a car by Blackmore  
Sensors and Analytics



Optical output gratings in silicon

# Silicon photonics solution

[www.sandia.gov/mesa/nspc](http://www.sandia.gov/mesa/nspc)  
[photonics@sandia.gov](mailto:photonics@sandia.gov)

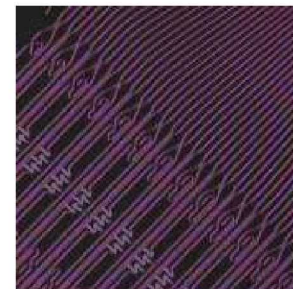
- Advantages

- Low-cost:** leverage investment on CMOS electronics
- Reliable:** no mechanical moving parts
- Compact:** several chips to provide large coverage
- Mature:** many device and system demonstrations

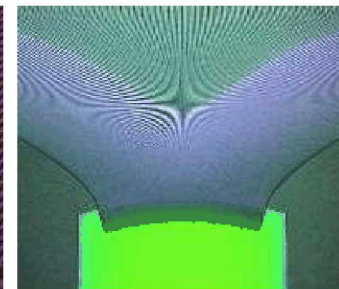
- Device challenges

- Thermo-optical-electronic packaging
- High optical power handling
- Fast optical phase error compensation
- Integration of new materials and layers

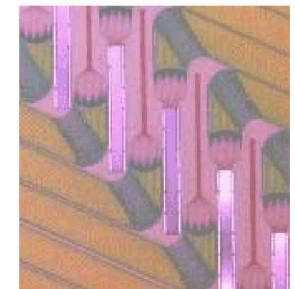
## Key components of a future silicon photonic LIDAR sensor



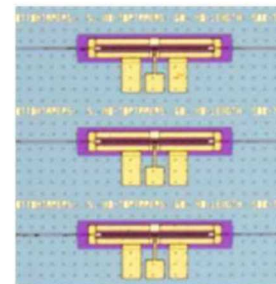
Low-loss, high-density waveguides



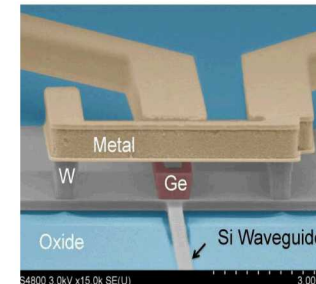
High-radix breakout



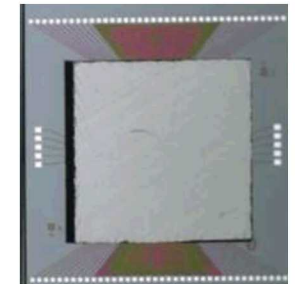
Efficient output couplers



Integrated laser and amplifiers



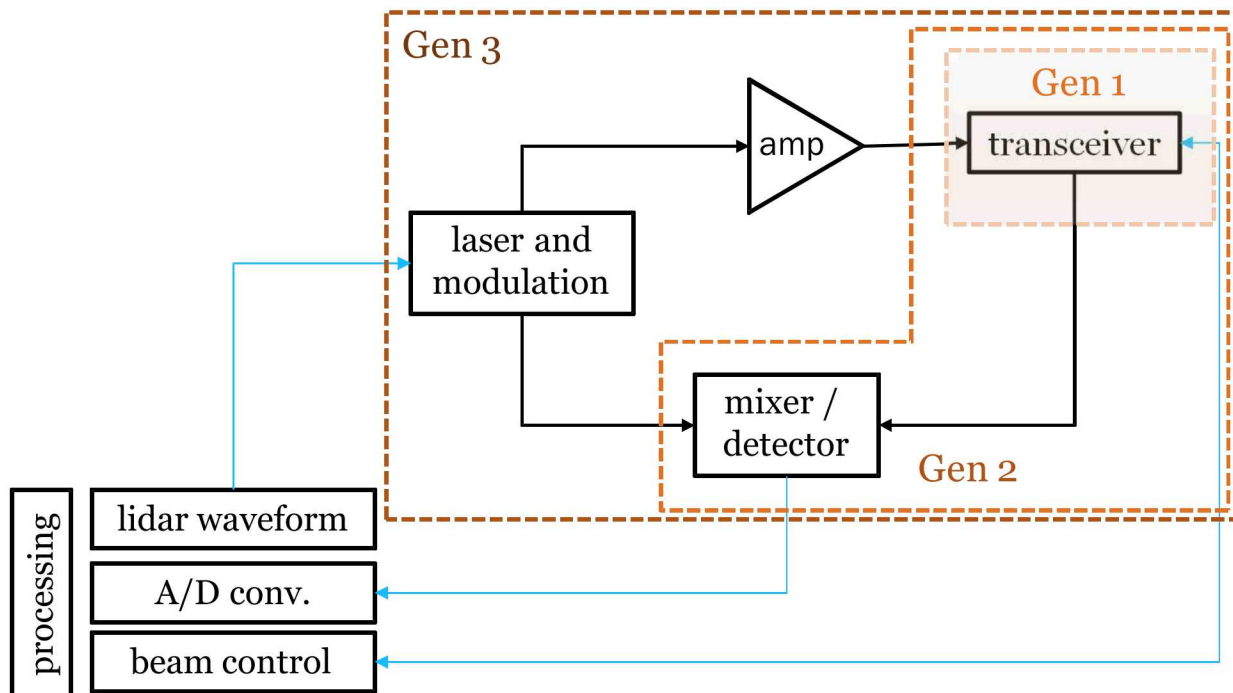
Ge photodetectors



Heterogeneous integrated CMOS electronics



# LiDAR engine



- **Gen 0 and 1**

- DARPA SWEEPER
  - proof of principle with beam scanners and a few emitters
  - patented idea for simplified controls
- Blackmore CRADA
  - expanded array size and added packaging
  - designs compatible with short-distance ranging

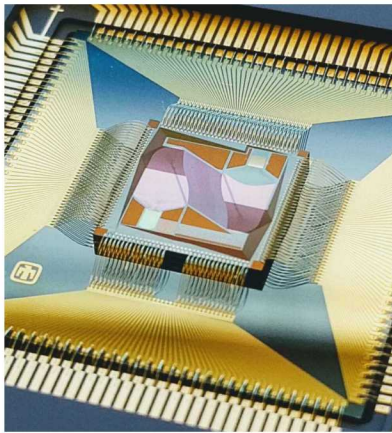
- **Gen 2**

- improved efficiency, power handling and functionality
- supported through LDRD program
- fabrication is underway at MESA fab

- **Gen 3**

- closer integration of laser, modulator, amplifiers, scanner, detector and CMOS controls to create a highly complex chip
- fab plan developed and patented

# Gen-1 results



Packaged 2D beam scanner

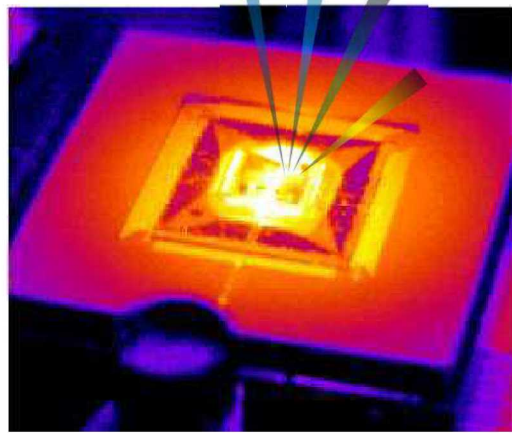
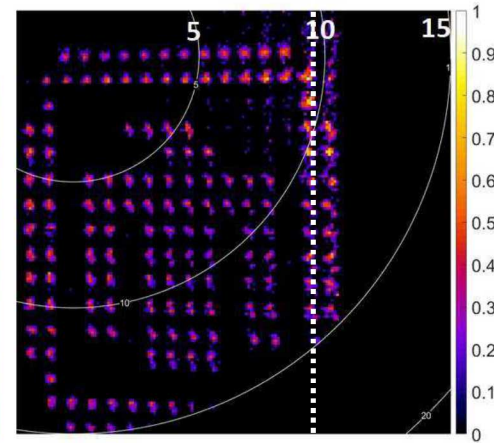
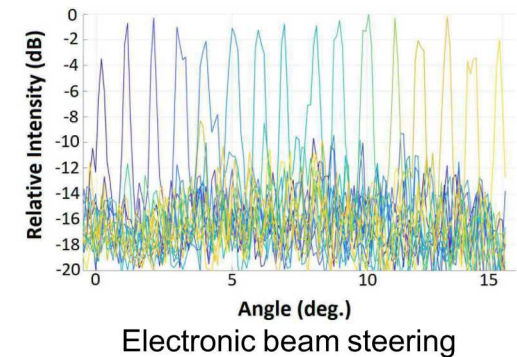


Image of chip with IR optical input



Composite image of beam scan

- 2D scanning with electronic and wavelength controls
  - electronic packaging with interposer and chip carrier
  - $N=256$  independent channels,  $d=3\text{ }\mu\text{m}$
  - long passive grating outputs for high fill factor aperture
- Field of view:  $24^\circ \times 10^\circ$ ; divergence angle:  $0.3^\circ \times 0.3^\circ$ 
  - near diffraction limited operation!
- Supporting technologies for new applications: machine vision, situational awareness, optical communication



M. Gehl, et al., CLEO 2019.



Exceptional service in the national interest



## Artificial Intelligence and Neural-Inspired Computing At Sandia National Laboratories

Presented by: William M Severa, PhD.

[wmsevera@sandia.gov](mailto:wmsevera@sandia.gov)



Sandia National Laboratories is a multi-mission laboratory managed and operated by National Technology and Engineering Solutions of Sandia LLC, a wholly owned subsidiary of Honeywell International Inc. for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.



## SANDIA HAS FIVE MAJOR PROGRAM PORTFOLIOS



Why has ML/AI had so much attention?

Some problems are difficult to solve with a directly-coded algorithm

- Don't generalize well
- Can be difficult to scale
- Have to write a program by hand for each specific task
  - Some tasks can be very difficult to encode
- Hand coded algorithms may run much slower

There have been Machine Learning (ML) successes in a variety of areas

- Recognizing patterns
- Anomaly detection
- Learning predictive models from data
- Creating surrogate models
- Automating repetitive computing tasks
- Generating synthetic data that models real data
- Assisting human decision making

These successes have been enabled by

- Large curated (labeled) datasets
- Advancements in computing power

### Dennard scaling

- As transistors get smaller, their power density remains constant

**Unfortunately ended 10-15 years ago**

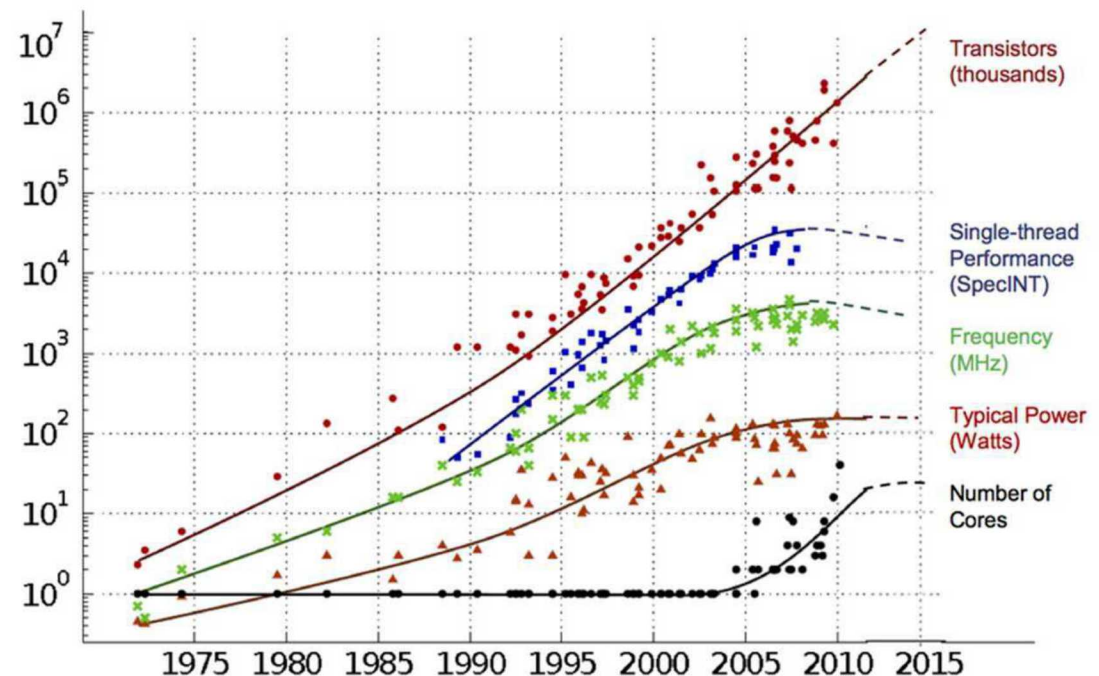
- Cannot run CPUs at faster speeds
- Emphasis on multi-core

Need for new paradigm of computing:

Novel Algorithms – Use AI to Accelerate

Novel Architectures – Accelerate AI

Novel Devices – Accelerate AI



Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten  
Dotted line extrapolations by C. Moore



# Sandia's Unique Mission Needs

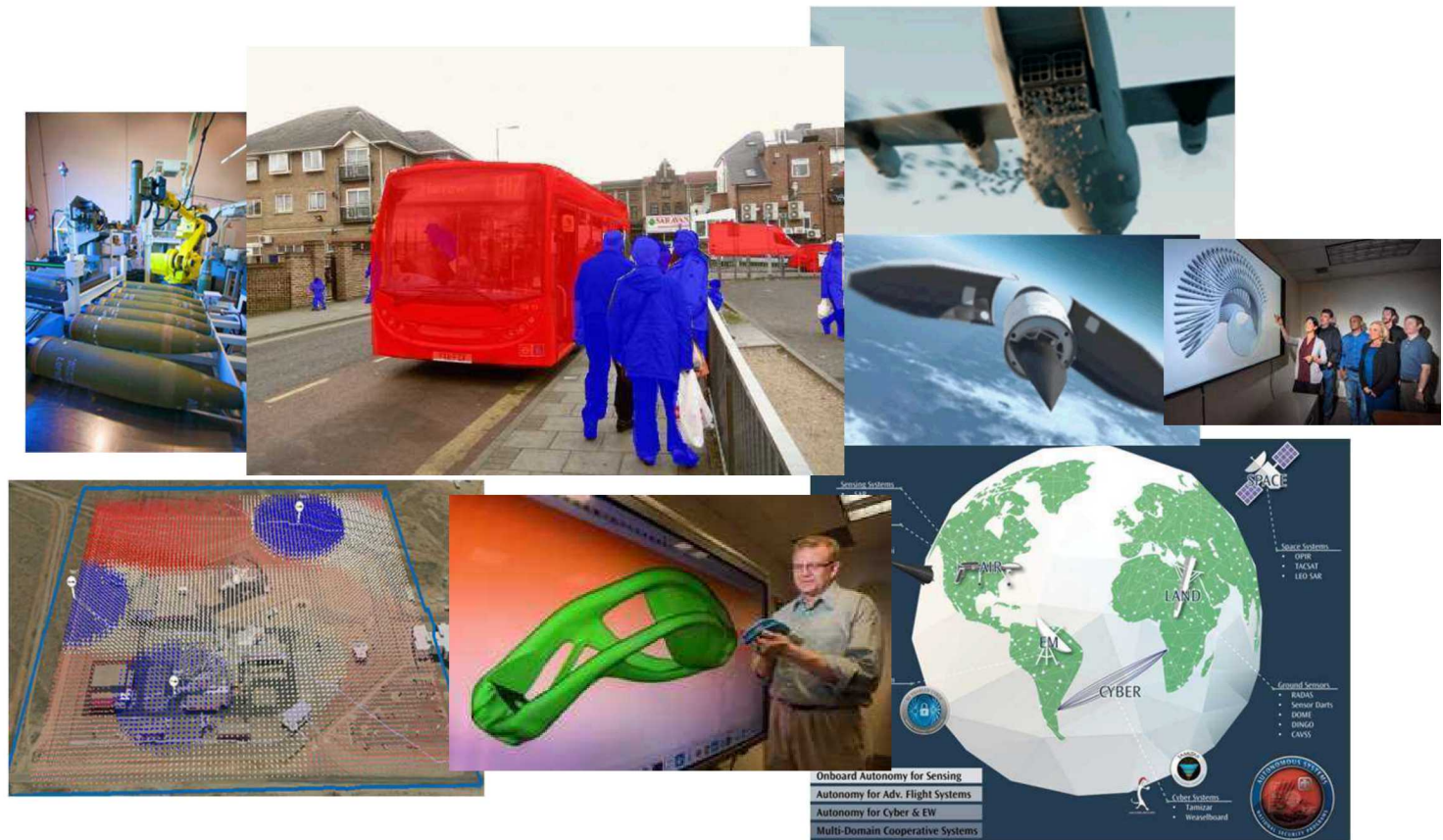


30

*Day*

*Scale*

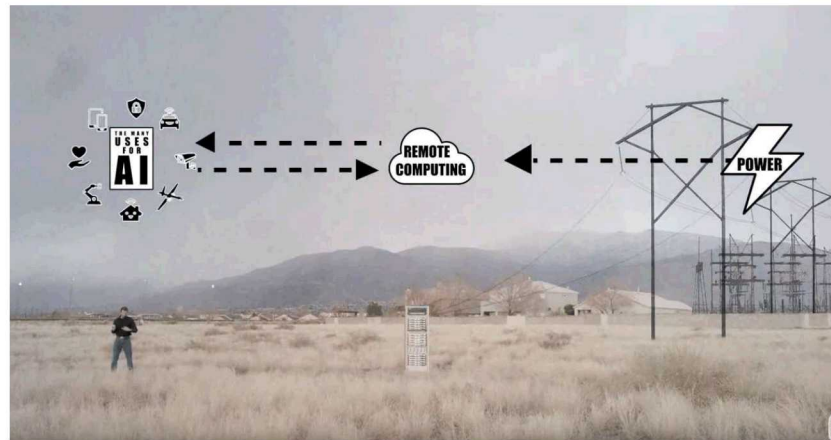
*Consequence*

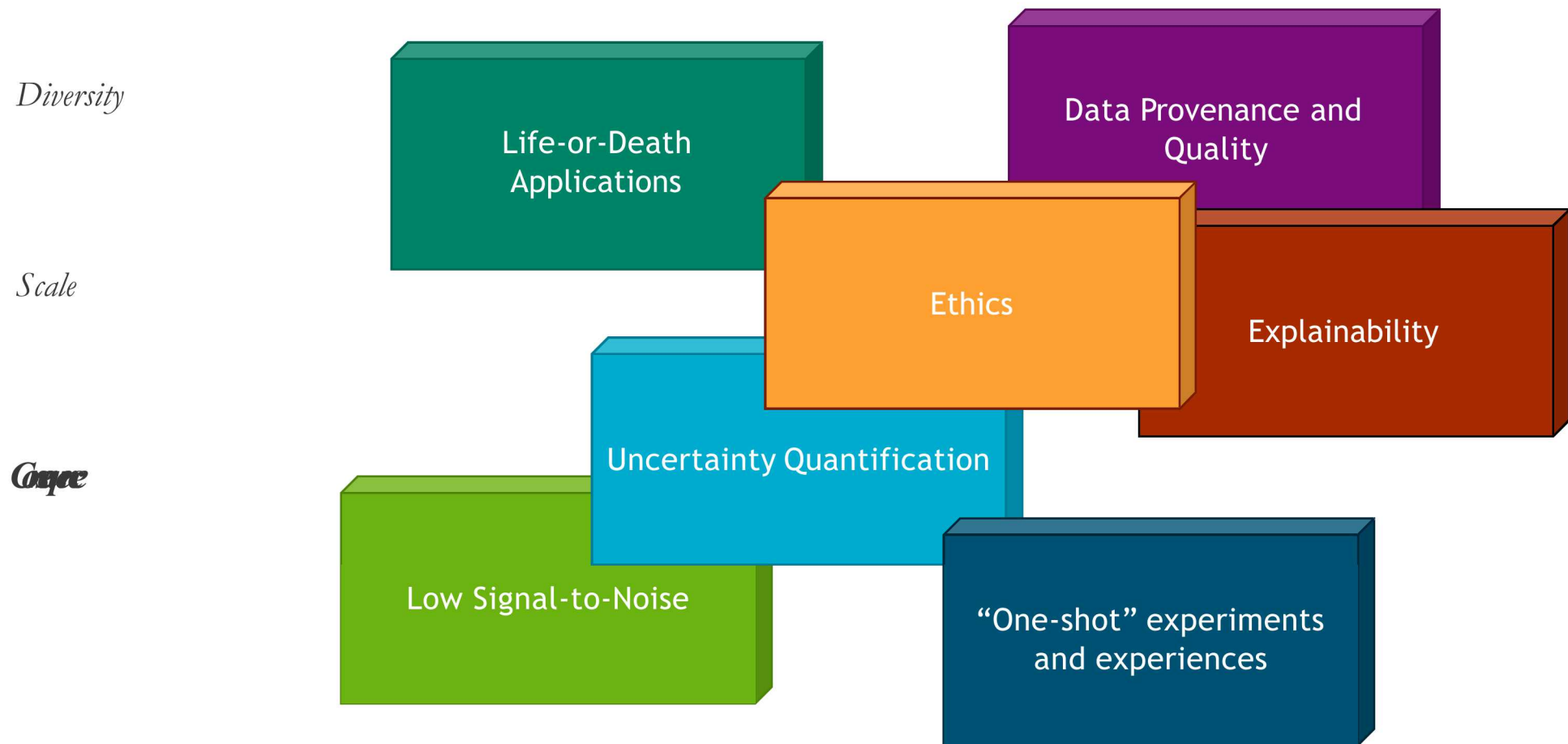


*Diversity*

*Sci*

*Consequence*







### High-confidence decisions

- Typically designing to “Five 9’s” of reliability
- Need to assure trust in our solutions
- Need to understand uncertainty of decisions
- Algorithms need to be explainable



 classified as rifle  
 classified as other

Synthesizing Robust  
Adversarial Examples,  
Athalye, et.al., 2018



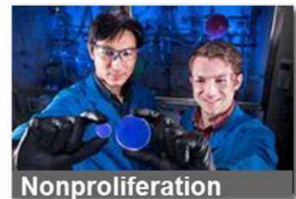
Military Systems



Space



Weaponneering



Nonproliferation



Infrastructure Resilience



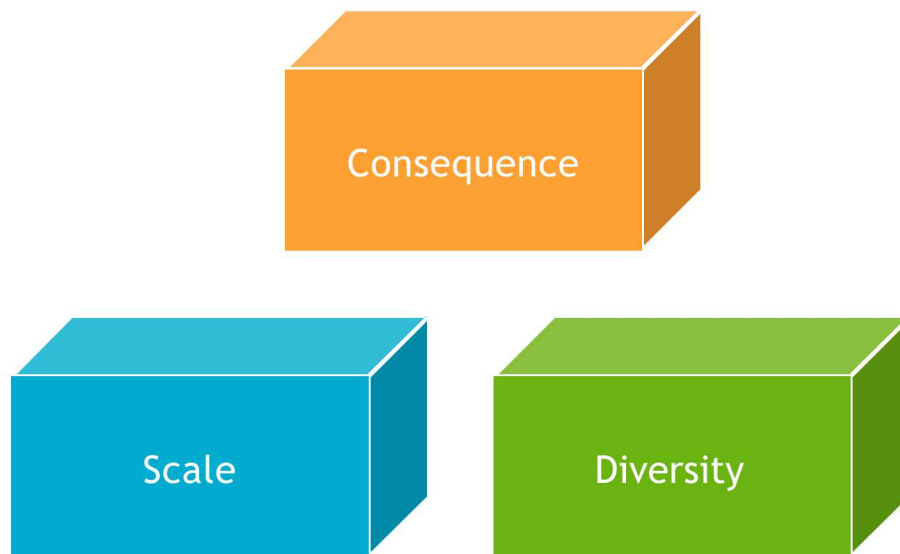
Homeland Security



Research

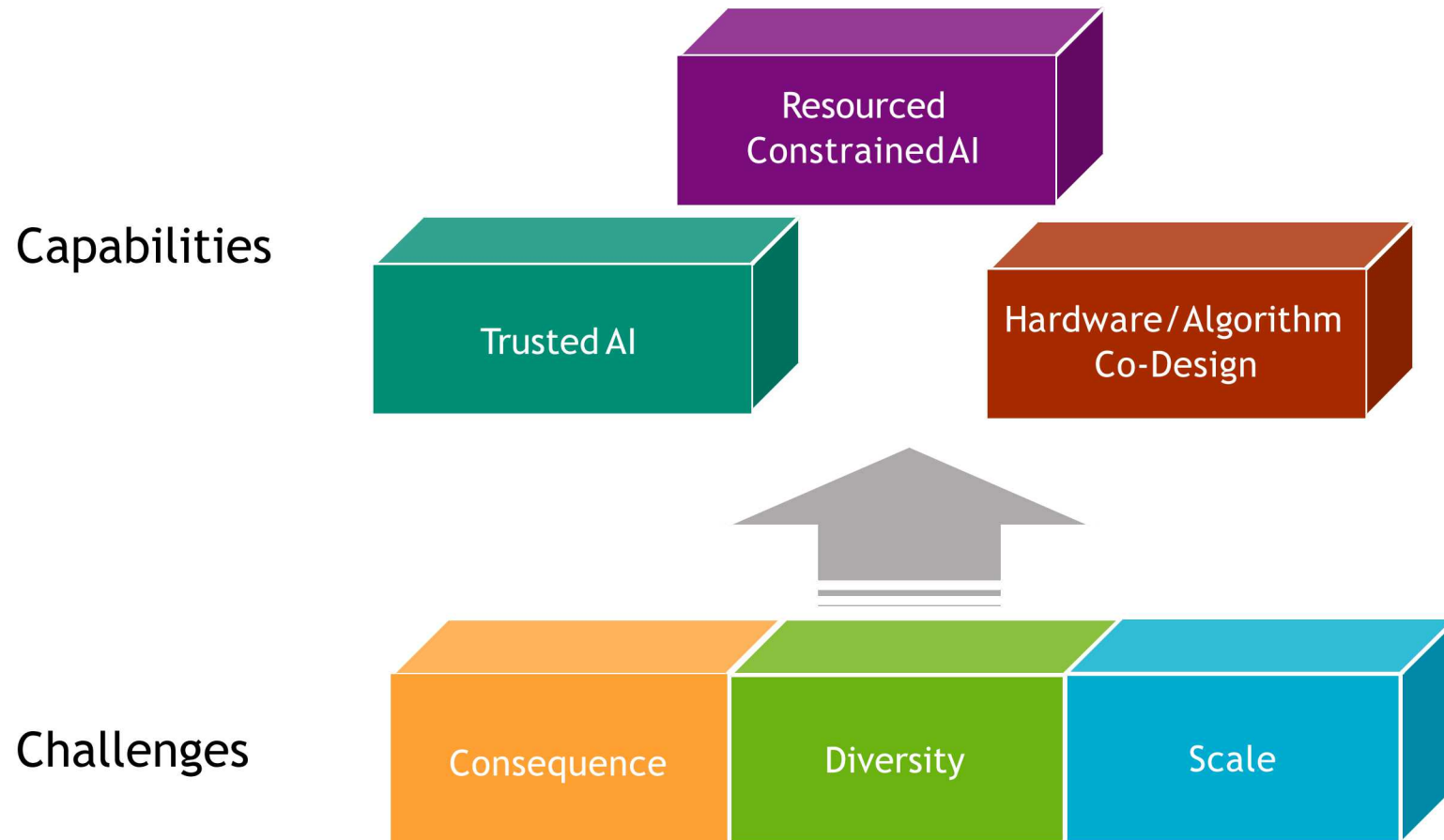


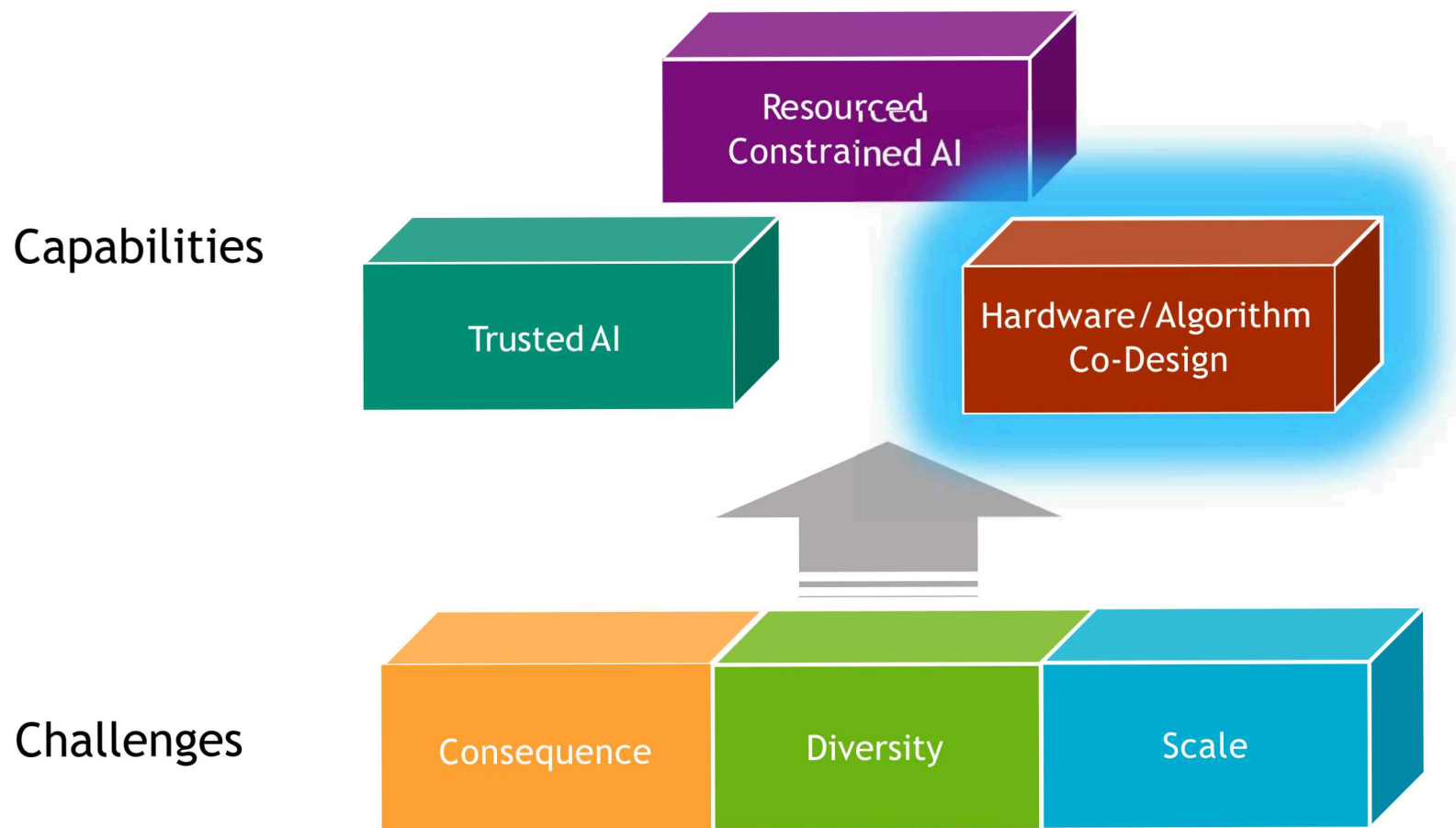
Many Sandia efforts are premised on idea that *brain-inspired* AI solutions will be instrumental in delivering these requirements



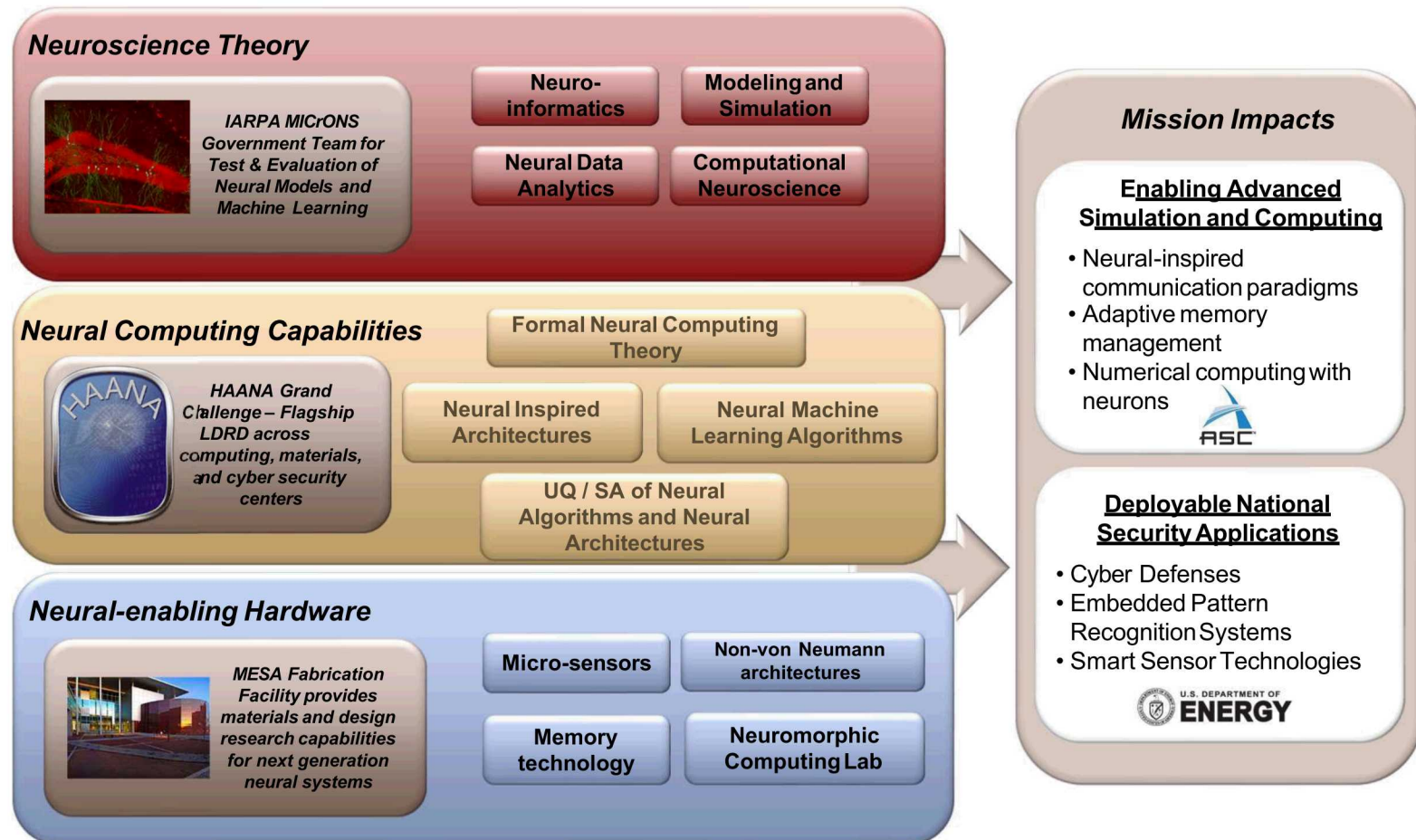
Sandia has a goal of creating a bridge between the broader world of AI and our missions

- Extending and developing AI algorithms
- Evaluating novel hardware and accelerators
- Explore brain-inspired sensor technology
- Identifying opportunities for novel AI impact
- Developing tools and analyses suitable for widespread adoption of emerging AI technologies





# Neural Computing at Sandia Labs Leverages a Large Research Foundation





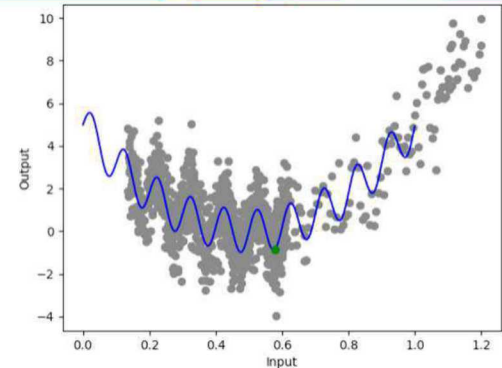
## AI/ML to Advance HPC Mission at Sandia

### Machine learning will provide new capabilities for scientific and engineering applications

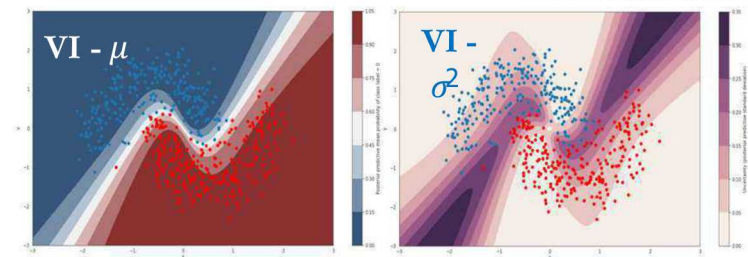
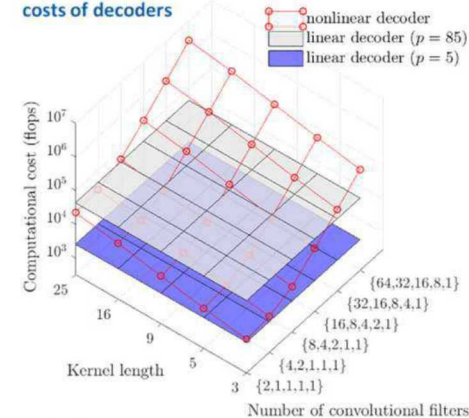
- Reduced order surrogate models for scientific/engineering problems
  - Could help us learn what is wrong/missing in physics models and aid in experimental design
- Ability to identify anomalies and regions of interest in inspection, surveillance, and large scale computational data
- Correlating and certifying simulation and experimental results
- Improving agility of application workflows (automating processes)

### Machine learning will provide new capabilities for HPC system administrators, facilities, dev-ops, and system software

- Help model complex behaviors (e.g., failures, degradation, energy)
- Automate/adapt usage to comply with more complex policy (e.g., energy consumption)
- Adaptable resource management (e.g., network, memory, storage, energy)
- “Smart” data-movement for Exascale runtimes



Analytically estimated computational costs of decoders



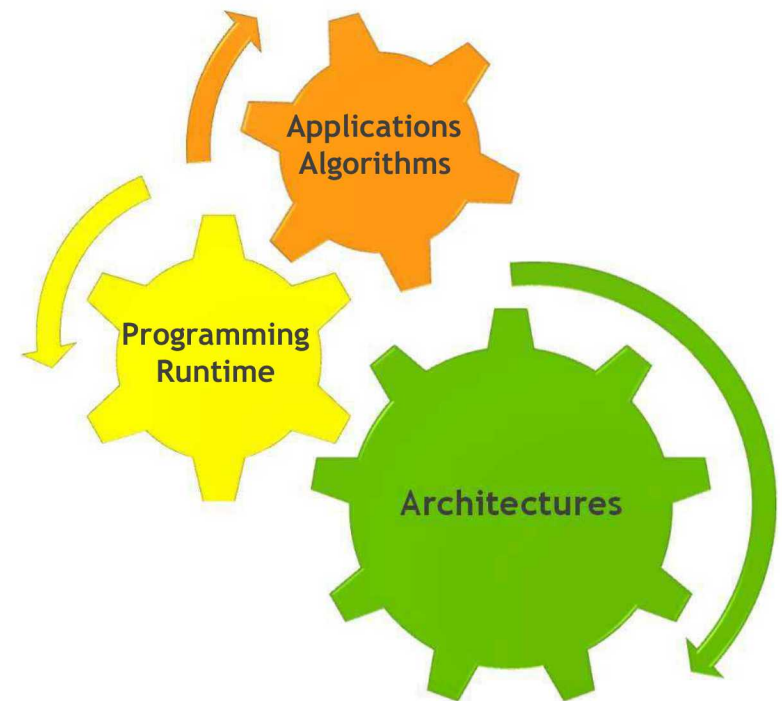
## Center for co-design of **AR**tificial Intelligence focused **AR**chitectures and **AR**gorithms (ARIAA)

ARIAA is a co-design research center that includes Pacific Northwest National Lab (PNNL), SNL, and Georgia Tech., supported by NVIDIA and Qualcomm

- Siva Rajamanickam, SNL PI (PNNL is lead lab)

ARIAA's objectives:

- Co-design novel AI/ML architectures, algorithms, and programming abstractions to enable traditional and ML-based DOE applications
- Understand how AI-focused dataflow/spatial architectures can impact future leadership class systems
- Understand how AI/ML accelerators can work with sparse, irregular, and/or streaming data



**Codesign of AI/ML accelerators with algorithms and applications will enable the development of this key technology to suit DOE HPC and AI/ML needs**



# Neuromorphic Hardware @ Neural Exploration & Research Lab (NERL)

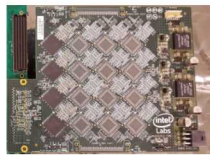


40

- ❑ Enables researchers to explore the boundaries of neural computation
- ❑ Consists of a variety of neuromorphic hardware & neural algorithms providing a testbed facility for comparative benchmarking and new architecture exploration



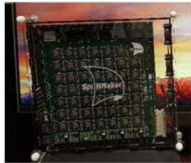
Intel Loihi



Intel Loihi



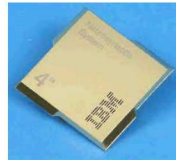
SpiNNaker 48 Node Board



SNL STPU on FPGA



IBM TrueNorth\*



Xilinx PYNQ FPGA



IBM TrueNorth NS16e\*



Nengo FPGA



Intel Neural Compute Stick



Nvidia Jetson TX1



\*Remote access

Google Coral



Nvidia Jetson Nano



Google EdgeTPU



GPU Workstations



Inilabs DAVIS 240C DVS



Cognimem CM1K



Georgia Tech FPAA



KnuPath Hermosa

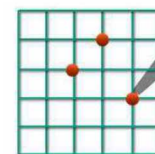




## Scientific Computing

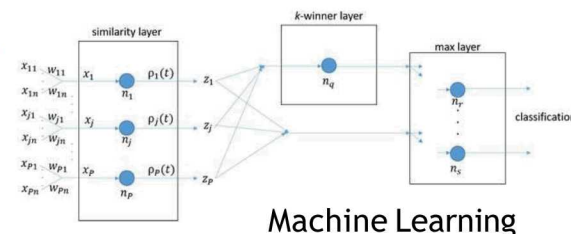
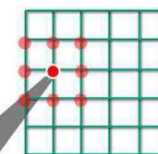
### Particle Method

Circuit per walker



### Density Method

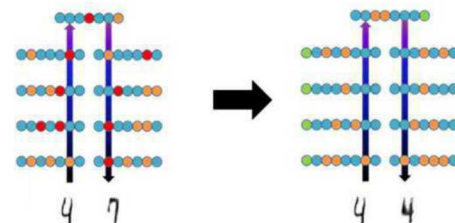
Circuit per position



## Machine Learning

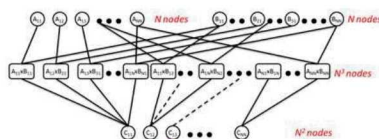


## Intelligent Storage

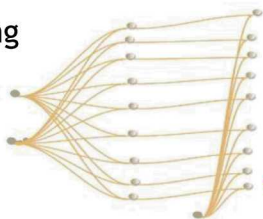


## Adaptive Deep Learning

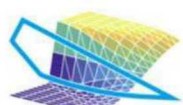
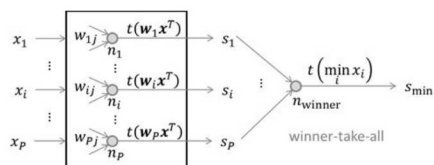
## Linear Algebra



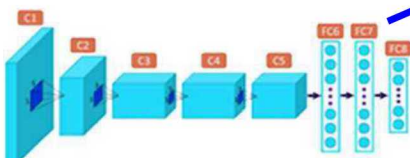
## Pattern Matching



## Optimizations



## WHETSTONE



## Context Modulated Deep Learning

## SNN

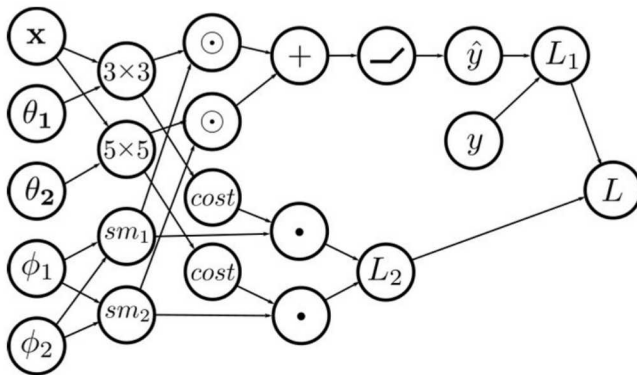
## Neural Algorithms

## NN

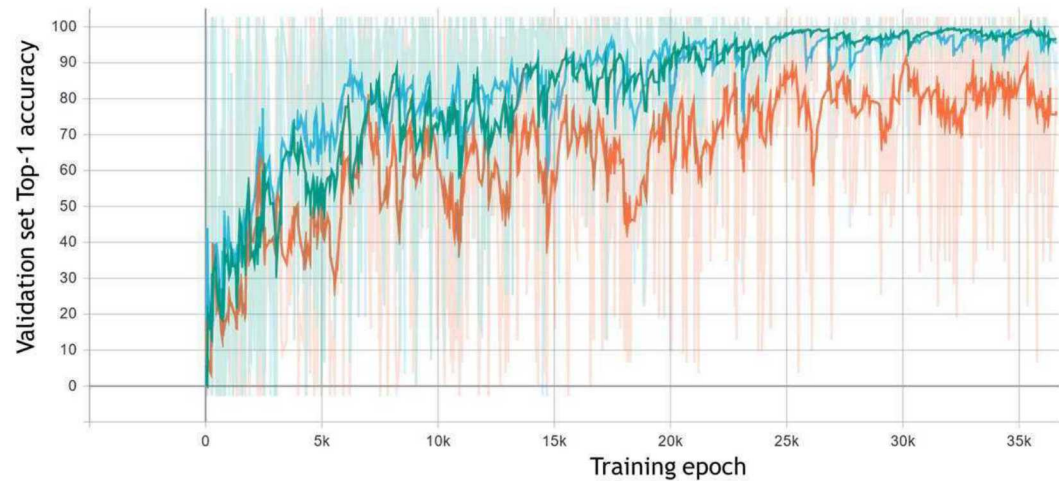
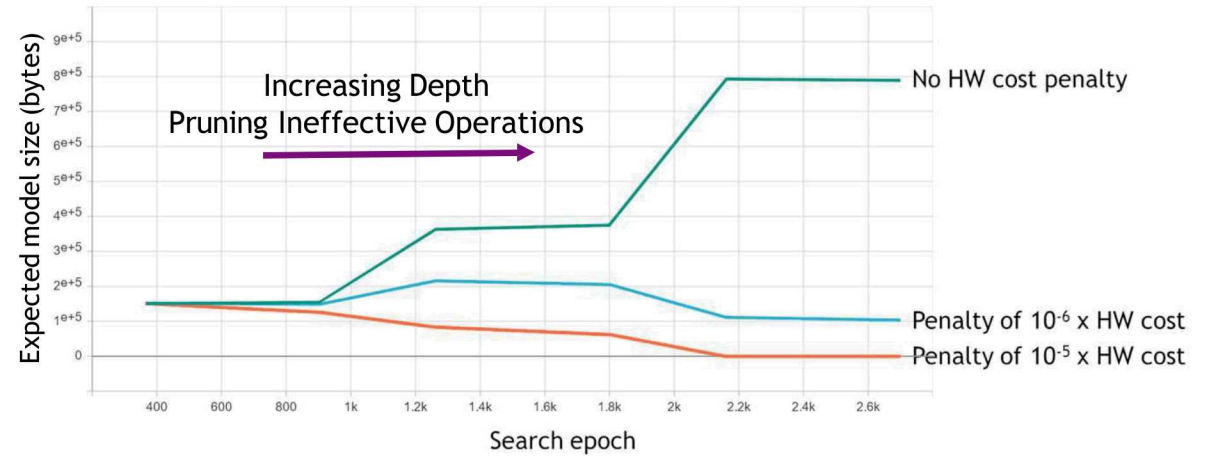
## ANN

# FORGE: Resource-Aware, Gradient-Assisted Neural Architecture Search

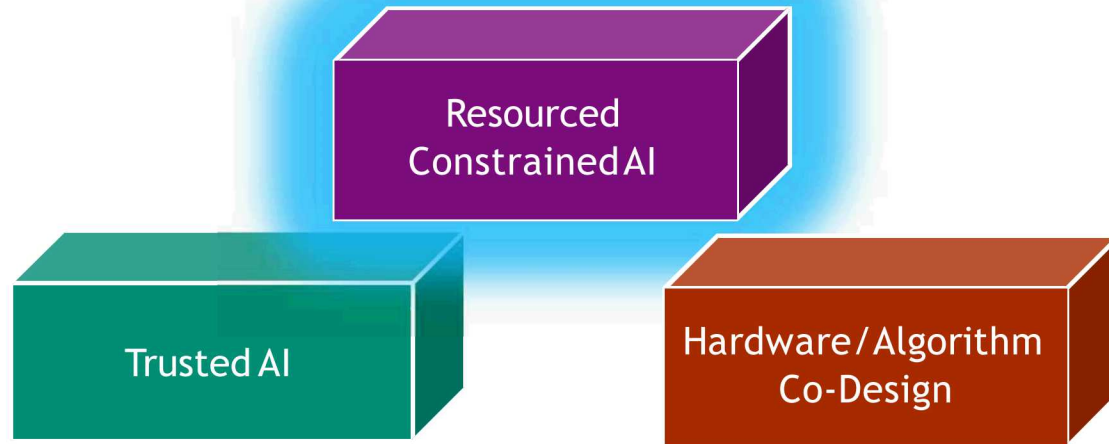
- Multi-level, Multi-objective Gradient-Assisted Optimization
- Automated Deep Learning Neural Network Design
- Highly Parallel Distributed Design
- Tailors Algorithm to Both Task and Target



Sample Network Schematic



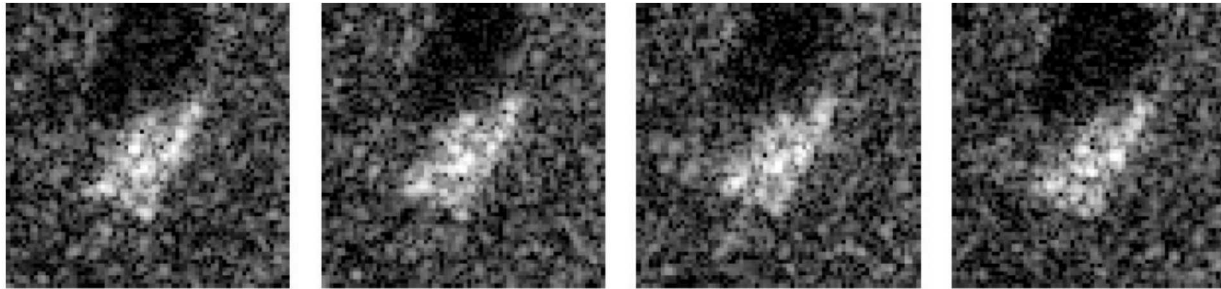
Capabilities



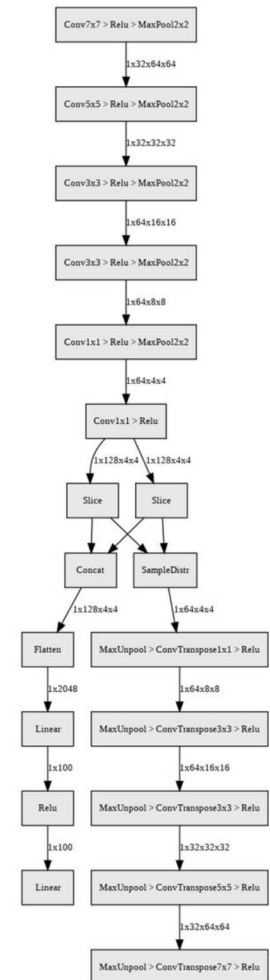
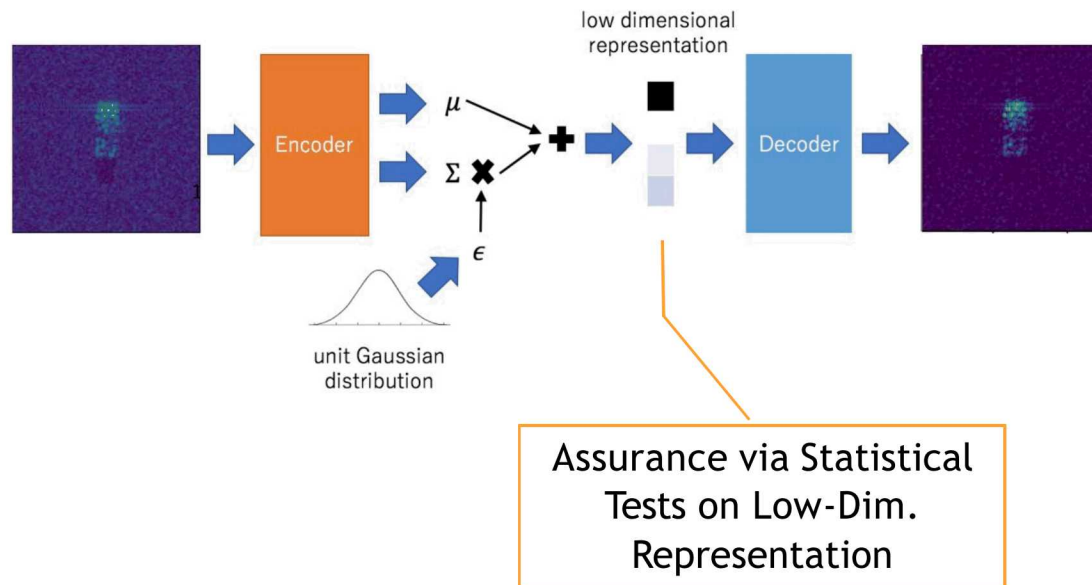
Challenges





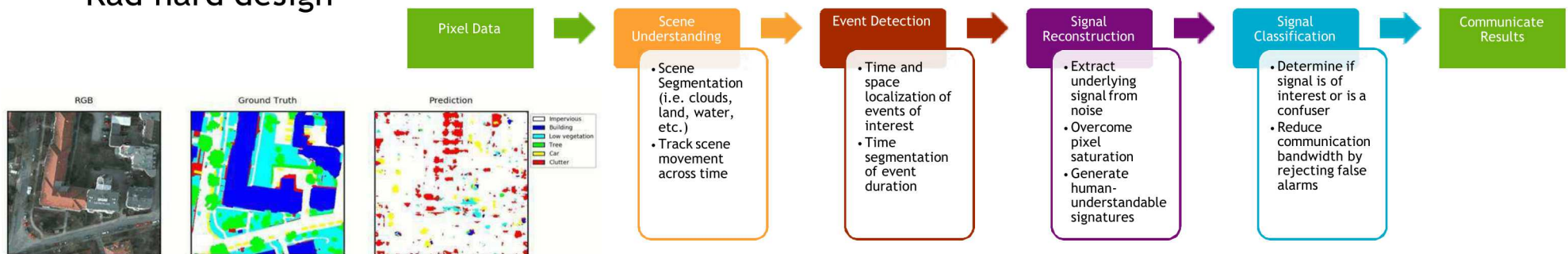
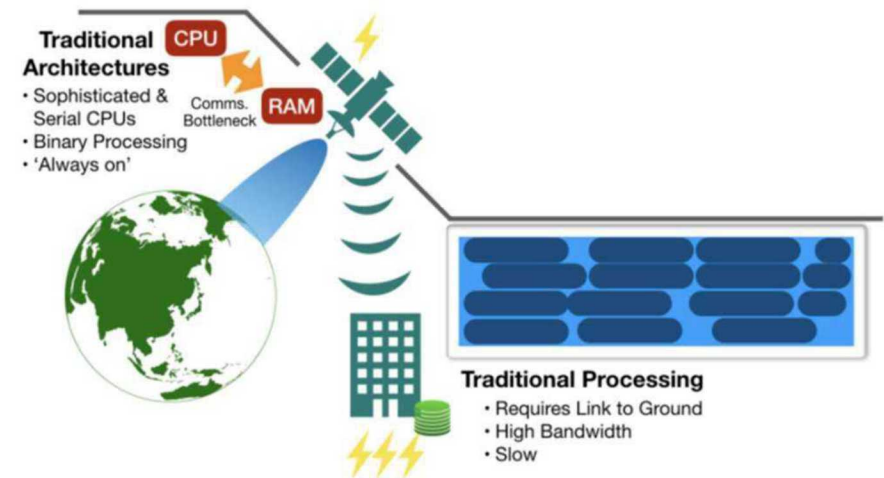


Synthetic Aperture Radar returns suffer signal variability due to coherence, specularity, and speckle



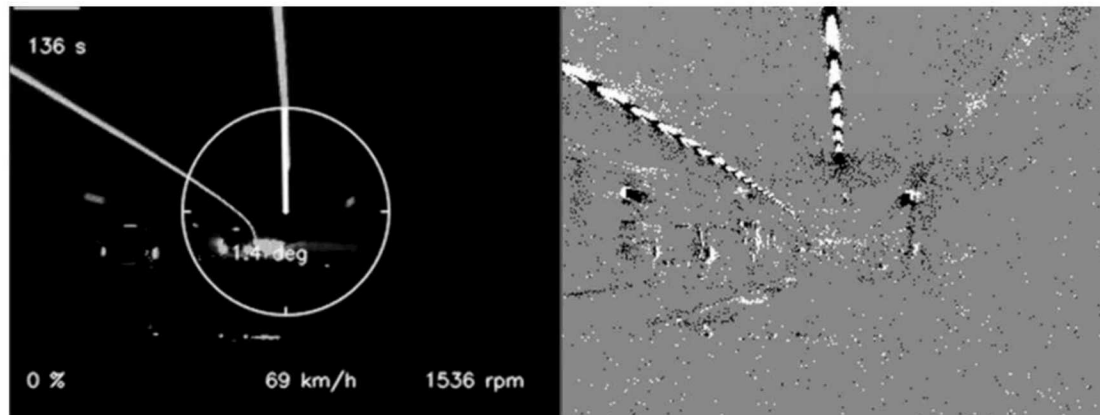
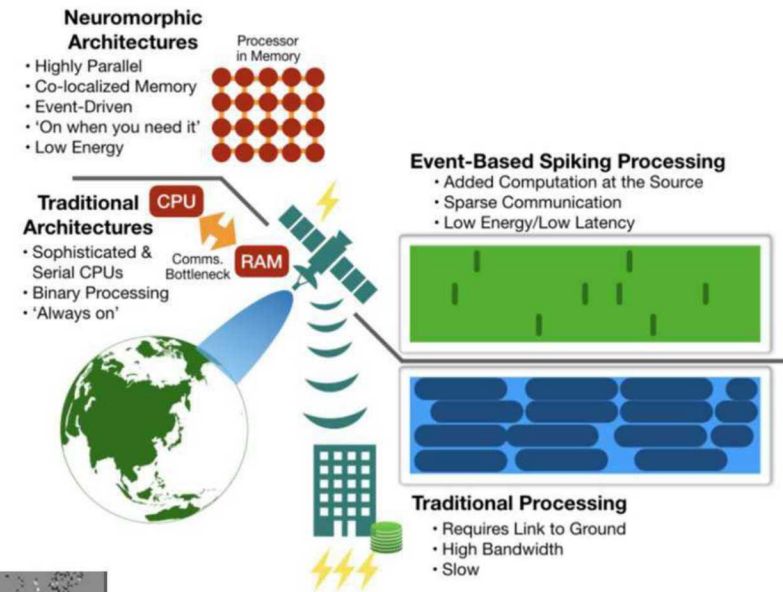
## Challenges in classic remote sensing

- Growth of sensor technologies outpacing communication bandwidth
- High Consequence Decisions
- Limited algorithm capabilities
- Limited onboard processing capability
- Rad hard design

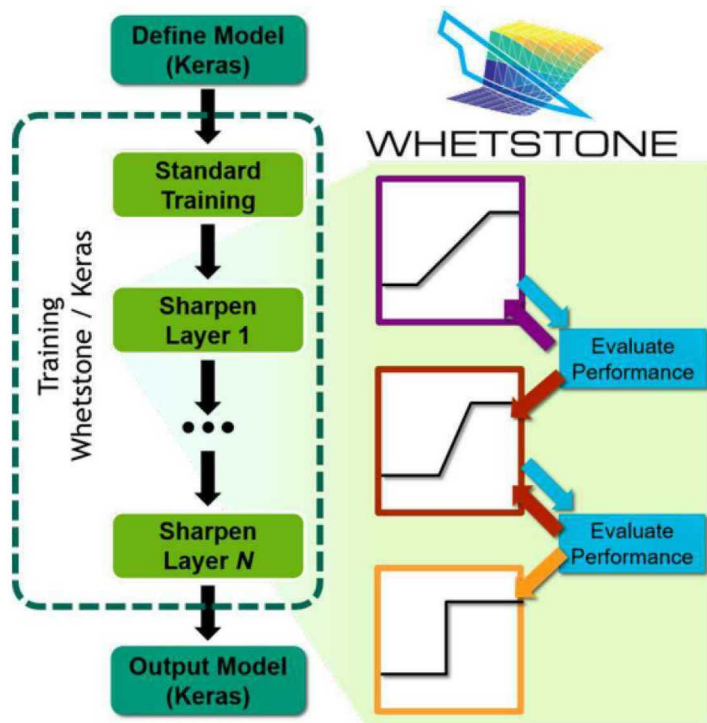


## Compute at the sensor

- Improve bandwidth utilization (send only what you need)
- Distributed computation avoiding single-point of failure
- May reduce preprocessing required (e.g. whitening)



DAVIS 240C  
Event-Driven  
Camera



Automatically converts deep learning networks from continuous valued neurons to binary activations, making them compatible with neuromorphic hardware

Open sourced

Published in February

Beginning to port onto neuromorphic platforms

SpiNNaker Results look great

ARTICLES  
<https://doi.org/10.1038/42756a>  
 nature machine intelligence

## Training deep neural networks for binary communication with the Whetstone method

William Severa<sup>1</sup>, Craig M. Vineyard<sup>2</sup>, Ryan Dellana<sup>3</sup>, Stephen J. Verzi<sup>4</sup> and James B. Aïme<sup>5</sup>\*

The computational cost of deep neural networks presents challenges to broadly deploying these algorithms. Low-power and embedded neuromorphic processors offer potentially dramatic performance-per-watt improvements over traditional processors. However, programming these brain-inspired platforms generally requires platform-specific expertise. It is therefore difficult to achieve state-of-the-art performance on these platforms, limiting their applicability. Here we present Whetstone, a method to bridge this gap by converting deep neural networks to have discrete, binary communication. During the training process, the activation function at each layer is progressively sharpened towards a threshold activation, with limited loss in performance. Whetstone sharpened networks do not require a rule code or other spike-based coding scheme, thus producing networks comparable in timing and size to conventional artificial neural networks. We demonstrate Whetstone on a number of architectures and tasks such as image classification, subvocoders and semantic segmentation. Whetstone is currently implemented within the Keras wrapper for TensorFlow and is widely extensible.

Artificial neural network (ANN) algorithms, specifically deep convolutional networks (DCNs) and other deep learning methods, have become the state-of-the-art techniques for a number of machine learning applications<sup>1–3</sup>. While deep learning models can be expensive both in time and energy to operate and even more expensive to train, their exceptional accuracy on fundamental analysis tasks such as image classification and audio processing has made their use essential in many domains.

Some applications can rely on remote servers to perform deep learning calculations; however, for many applications, such as onboard processing in autonomous platforms like self-driving cars, drones and smart phones, the resource requirements of training large ANNs may still prove to be prohibitive<sup>4,5</sup>. Large ANNs with many parameters require a significant storage capacity that is not always available, and data movement energy costs are greater than that of performing the computation, making large ANNs intractable<sup>6</sup>. Additionally, onboard processing capabilities are often limited to meet energy budget requirements, further complicating the challenge. Other factors such as privacy and data sharing also provide a motivation for performing computation locally rather than on a remote server.

The development of specialized hardware to enable more efficient ANN calculations seeks to facilitate moving ANNs into resource-constrained environments, particularly for trained algorithms that simply require the deployment of an inference-ready network. A common approach today is to optimize key computational kernels of ANNs in application-specific integrated circuits (ASICs)<sup>7–9</sup>. However, while these ASICs can provide substantial acceleration, their power costs are still too high for some embedded applications and often lack flexibility for implementing alternative ANN architectures.

Brain-inspired neuromorphic hardware presents an alternative to conventional ASIC accelerators, and has been shown to be capable of running ANNs with potentially orders-of-magnitude lower power consumption (that is, performance-per-watt). The landscape of neuromorphic hardware is rapidly evolving<sup>10–12</sup>; however, increasingly these approaches leverage spiking to achieve substantial energy

savings. Neuromorphic spiking, which emulates all-or-none action potentials in biological neurons, limits communication in hardware only to discrete events. For spiking neuromorphic hardware to be useful, however, it is necessary to convert an ANN, for which communication between artificial neurons can be high-precision, to a spiking neural network (SNN). Supplementary Note 1 provides further details of spiking and ANN acceleration.

The conversion of ANNs to SNNs—whatever their form—is non-trivial, as ANNs depend on gradient-based backpropagation training algorithms, which require high-precision communication, and the resultant networks effectively assume the persistence of that precision. While there are methods for converting existing ANNs to SNNs, these transformations often require using representations that diminish the benefits of spiking. Here, we describe a new approach to training SNNs, where the ANN training is to not only learn the task, but to produce a SNN in the process. Specifically, if the training procedure can include the eventual objective of low-precision communication between nodes, the training process of a SNN can be nearly as effective as a comparable ANN. This method, which we term Whetstone (Fig. 1) inspired by the tool to sharpen a dull knife, is intentionally agnostic to both the type of ANN being trained and the targeted neuromorphic hardware. Rather, the intent is to provide a straightforward interface for machine learning researchers to leverage the powerful capabilities of low-power neuromorphic hardware on a wide range of deep learning applications (see section ‘Implementation and software package details’).

**Results**  
 Whetstone method converts general ANNs to SNNs. The Whetstone algorithm operates by incorporating the conversion into binary activations directly into the training process. Because most techniques to train ANNs rely on stochastic gradient descent methods, it is necessary that the activations of neurons be differentiable during the training process. However, as networks become trained, the training process is able to incorporate additional constraints, such as targeting discrete communication between nodes. With this kind of the optimization target in

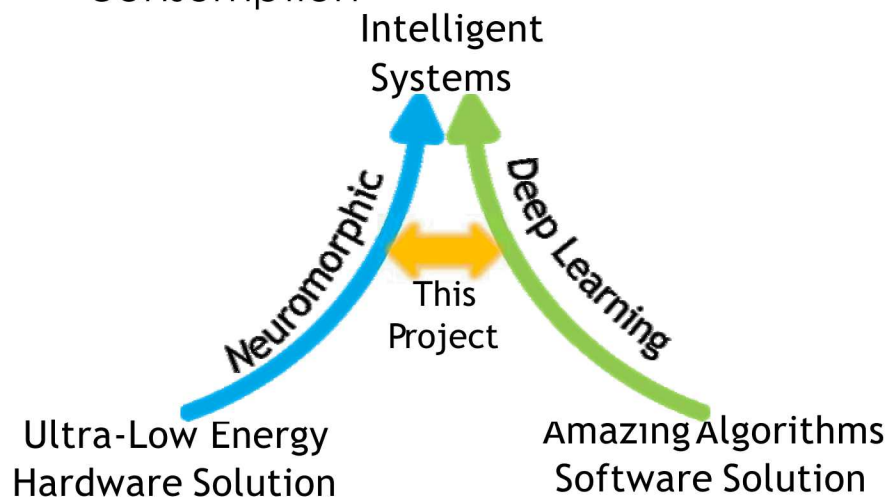
Center for Computing Research, Sandia National Laboratories, Albuquerque, NM, USA. \*e-mail: severa@cs.sandia.gov; jbaime@cs.sandia.gov

NATURE MACHINE INTELLIGENCE | VOL 1 | FEBRUARY 2019 | DOI:10.1038/s42256-019-0000-9

86



- AI power draw is a key limiting factor especially for electric powered vehicles: 3kW now; HPC-level for fully self-driving
- Prototype vehicles use a trunk full of GPUs
- Forecasting current tech ~1TeraOp/Watt
- Neuromorphic Hardware:
  - Enables event-driven computation
  - Opportunity for extremely low power consumption



SpiNNaker, Univ. of Manchester

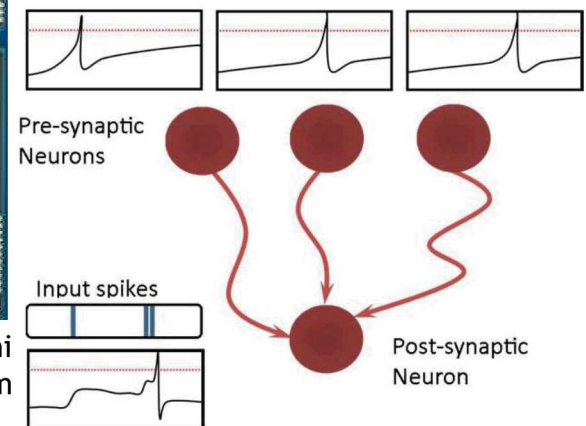


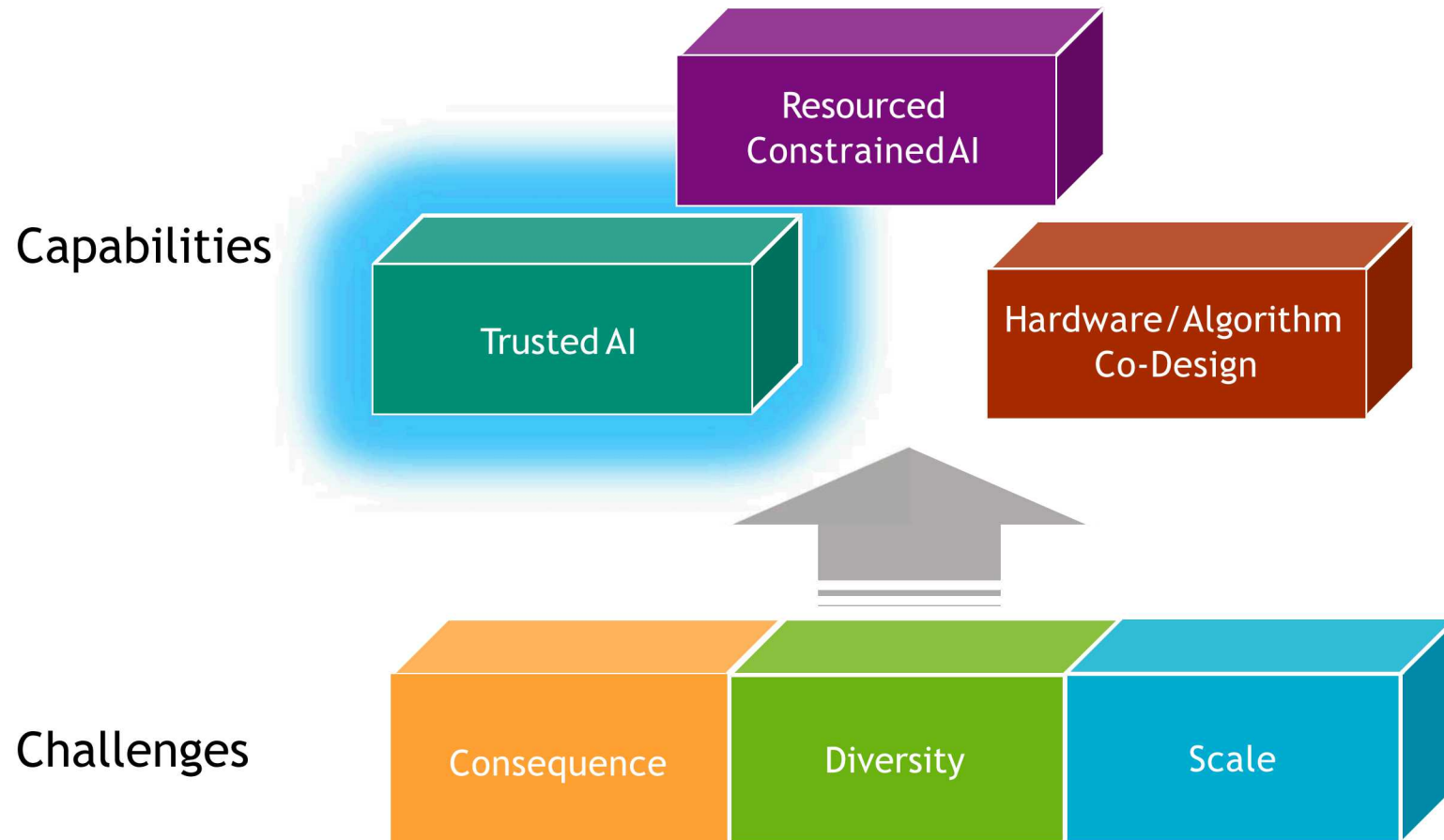
Example Image  
From Berkeley DeepDrive



Intel Loihi  
Photo: intel.com

## Spiking Neuron Representation





## Problem

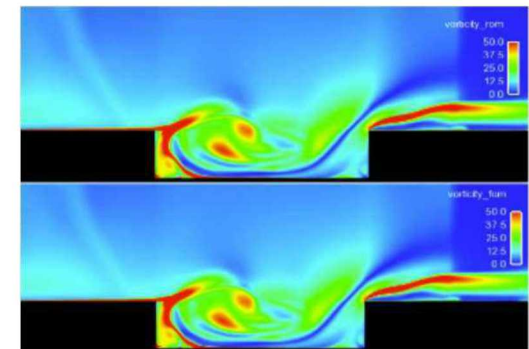
- High-fidelity computational physics simulations on HPC systems can take hours or days to execute
- Lengthy execution time limits the design space explored during conceptual design
- Need a faster, more efficient means of simulating complex physics problems

## Technical Approach

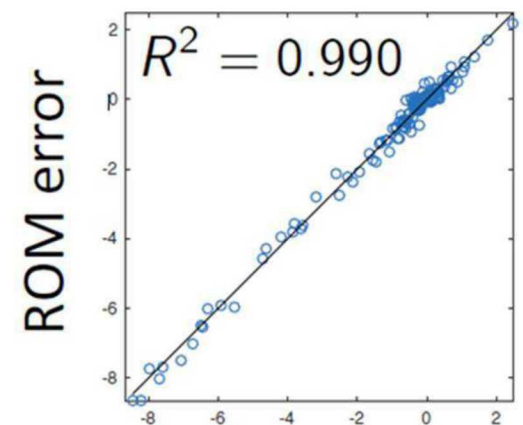
- Create Reduced Order Model (ROM) from high-fidelity simulation data that
  - Executes faster via dimensionality reduction using autoencoders without significant reduction in accuracy
  - Preserves important physical properties (e.g., conservation laws)
  - Uses Machine Learning Error Models (MLEM) to quantify uncertainty

## Results/Accomplishments

- Reduced order surrogate models and theory have been developed for turbulent flow simulations
- Runtimes are 100-1000 times faster and are only 1% less accurate than the high-fidelity simulations
- MLEM can predict errors with validated statistical properties



Turbulent flow vorticity field



## Problem

High-fidelity simulations on HPC systems are too expensive

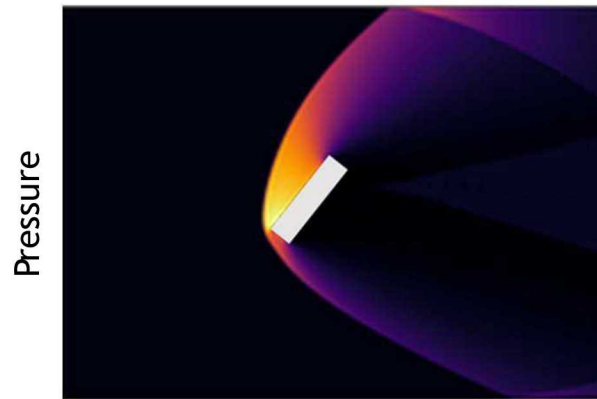
## Technical

Train a neural network to predict the steady-state flow field  
Guide the prediction with physical constraints (conservation laws) and aerodynamic forces (drag, lift, torque)

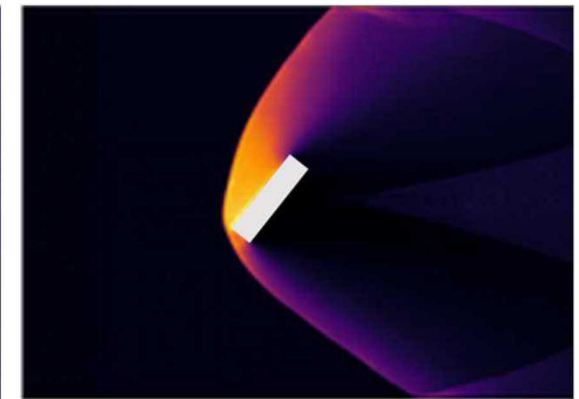
## Results/Accomplishments

Demonstrated >100x speed increase in 2D with < 6% average error  
Predict > 1000x speed increase in 3D

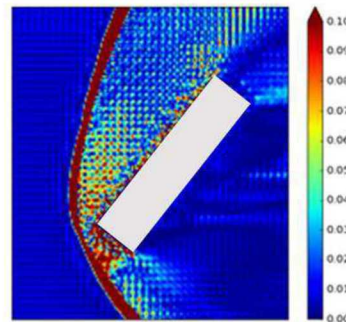
Hydro-code Simulation



ML Prediction



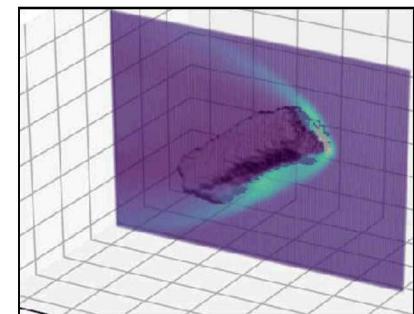
Pressure



Relative error map of ML prediction

2D force	Avg Error
Drag	1.87%
Lift	5.63%
Torque	2.29%

ML model successfully predicts flow field and aerodynamic coefficients

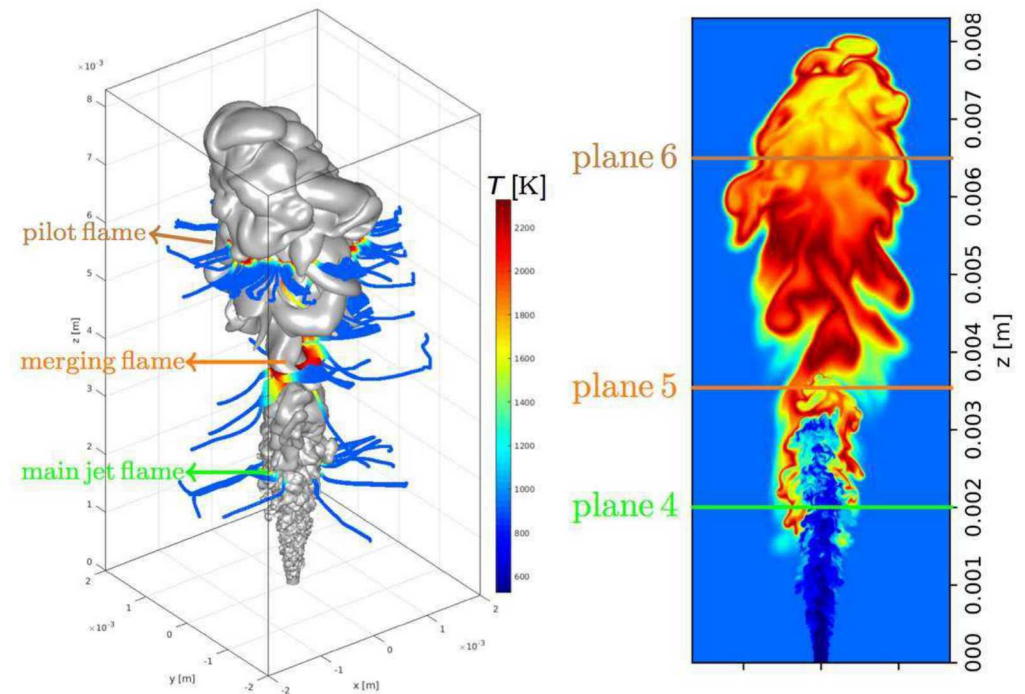


ML prediction of pressure field around complex 3D object



# Surrogate Reduced Order Modeling of Diesel Combustion and Ignition in GDI engines using ML/AI

- Surrogate reduced-order modeling with principal component transport of compositions to reduce the chemical dimensions needed to describe low- and high temperature ignition, flame propagation and soot under diesel and cold-start GDI conditions
- Use petascale DNS data and experiments to provide ‘truth’ data to adaptively train reduced-order surrogates incorporating physics constraints (e.g. using governing equations)
- Anomaly detection ML for detecting pre-ignition and knock
- Reduced order modeling for engine design and optimization



DNS of n-dodecane multi-injection diesel combustion showing instantaneous volume rendering of mixture fraction (left image) and temperature slice (right image). Multi-dimensional flamelets at 3 axial planes shown by blue lines (Rieth, Chen, Xu, Han, Hasse)

## Problem

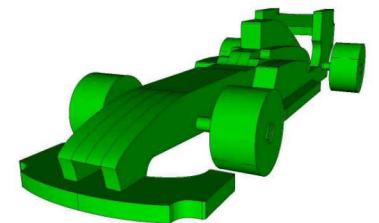
- Geometry preparation and meshing for computational simulation is bottleneck (consuming 70%+ of analyst time)
- Analyst/engineer must have extensive domain-specific expertise to manage many individual complex problems and tasks
- Must produce verifiably accurate physics appropriate mesh ready for simulation

## Technical Approach

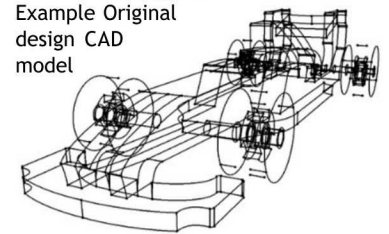
- Identify tasks currently done by analysts to train machine learning models
- Capture and label operations performed by expert using existing software
- Build a feature library of geometric characteristics commonly encountered in CAD models and identify solutions for effectively modifying CAD for best resulting mesh
- Explore machine learning models that provide best solutions for CAD features with associated solution labels

## Results/Accomplishments

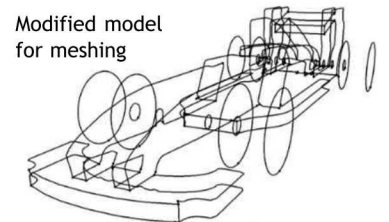
- Developed ML techniques to rank geometry-modification operations by their likelihood of yielding a meshable model
- Provides insight on which geometric features are most useful for machine learning, and would be relatively easy to integrate into the analyst workflow if successful



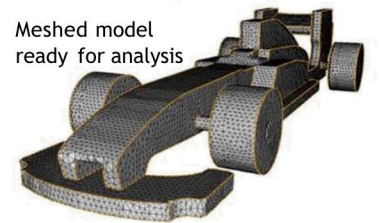
Example Original design CAD model



Modified model for meshing



Meshed model ready for analysis



## Problem

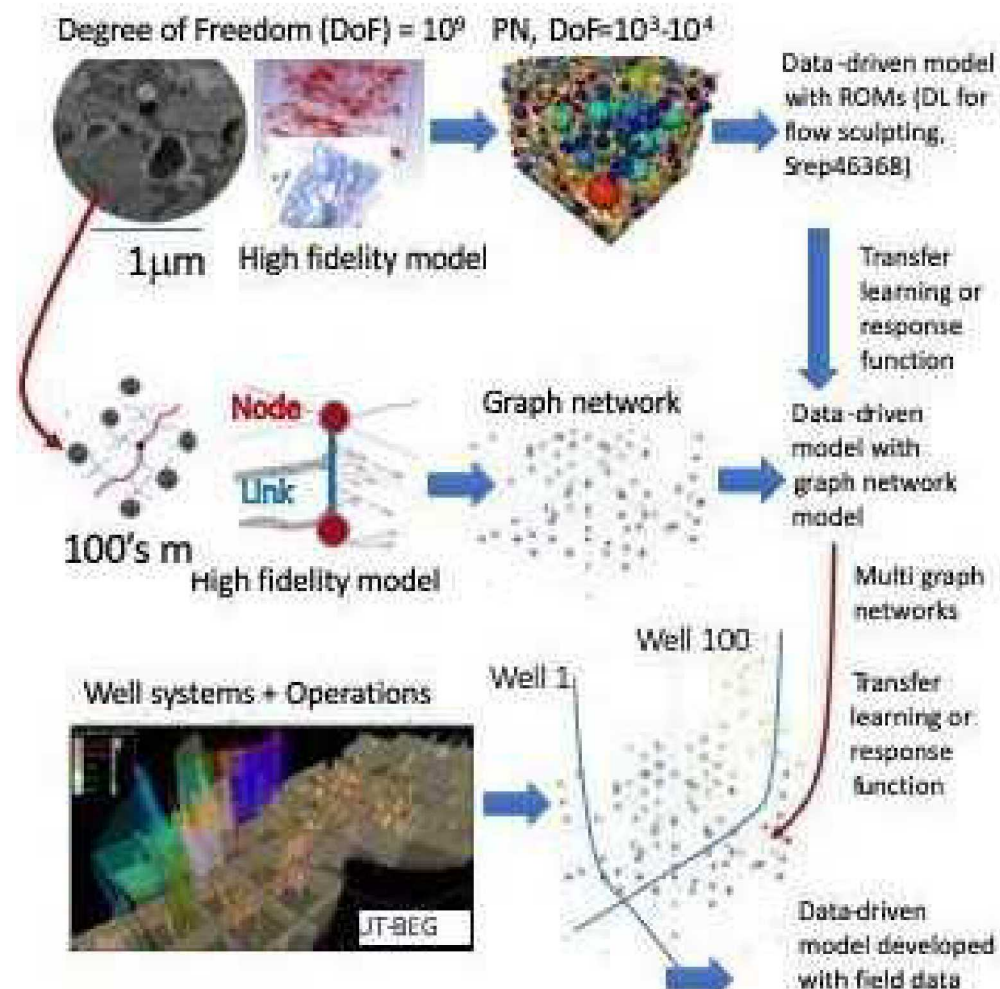
- Given sparse data at various scales (nm, m, km) about rocks and wells, and a given depth, can a machine learning algorithm predict how much shale gas can be extracted?

## Technical Approach

- Develop multiscale, physics informed deep learning algorithm to generate, parse, and predict data to solve the above problem
- Develop physics informed costs and constraints driven algorithms

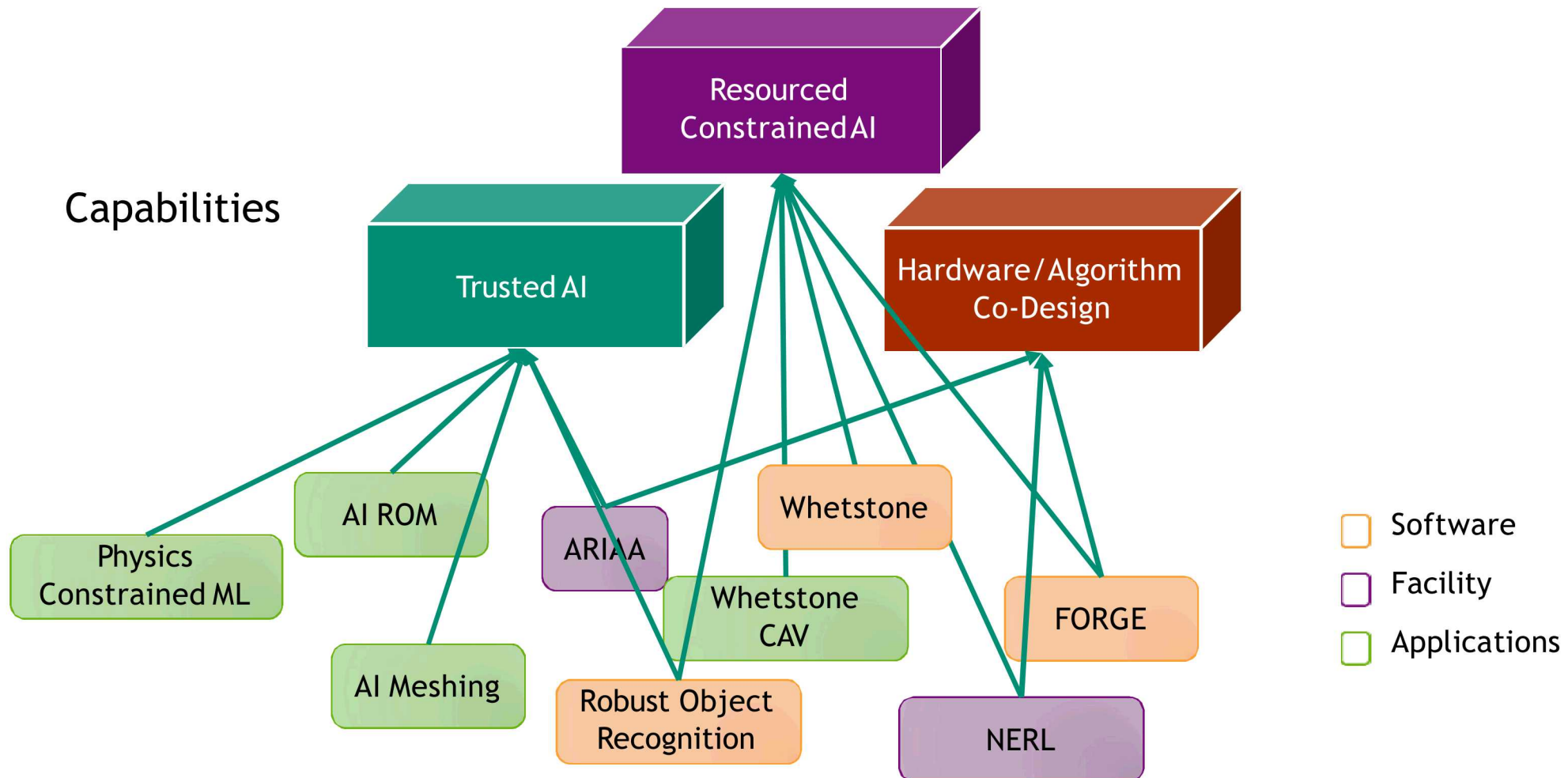
## Results/Impact

- Improve current state of art predictions and resource estimates
- Develop physics informed machine learning algorithms





## Capabilities





Thanks! Questions?



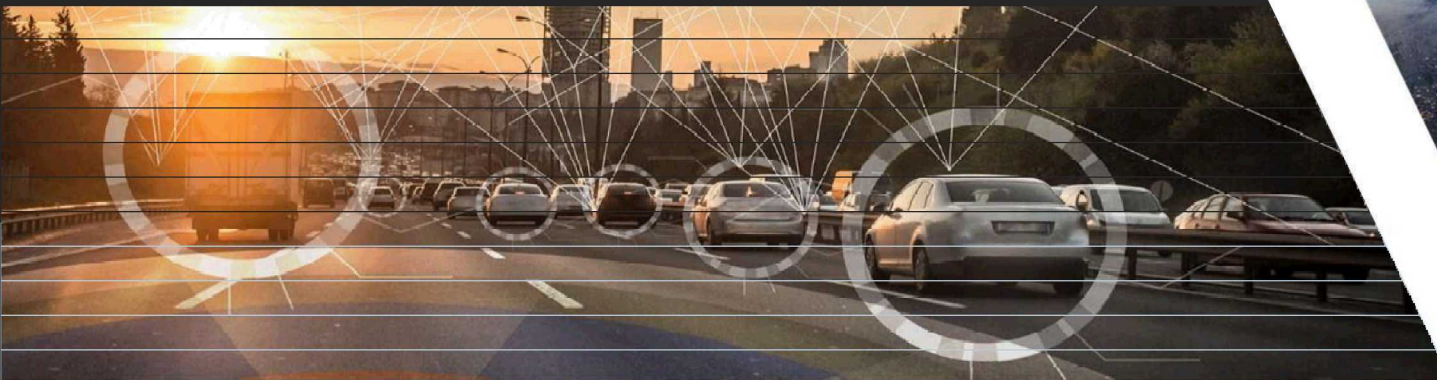
Exceptional service in the national interest





**Sandia  
National  
Laboratories**

*Exceptional service in the national interest*



## Advanced Computing for Connected & Automated Vehicle (CAV)

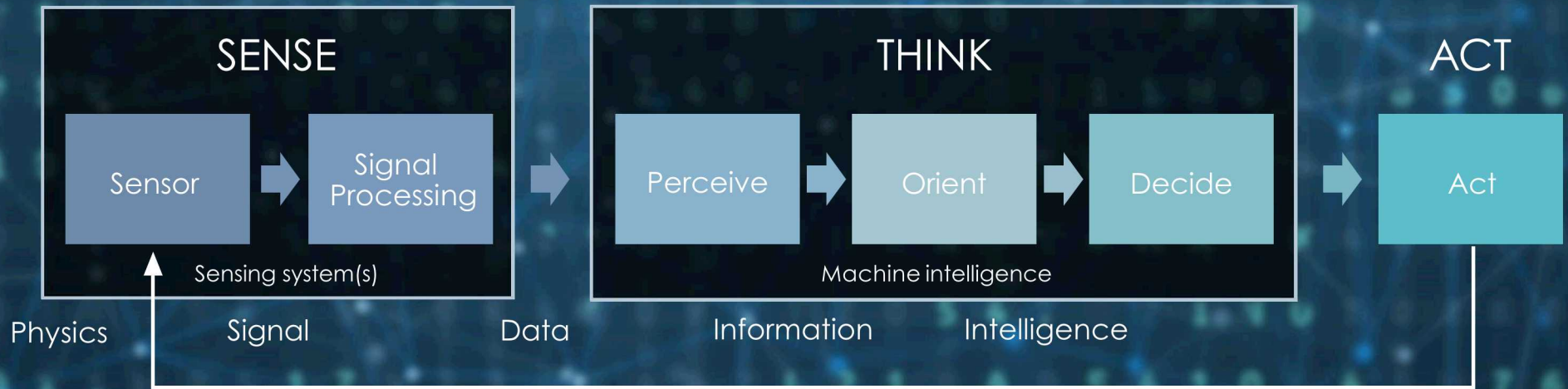


Sandia National Laboratories is a multission laboratory managed and operated by National Technology and Engineering Solutions of Sandia LLC, a wholly owned subsidiary of Honeywell International, Inc. for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.





# AUTOMATED SYSTEMS





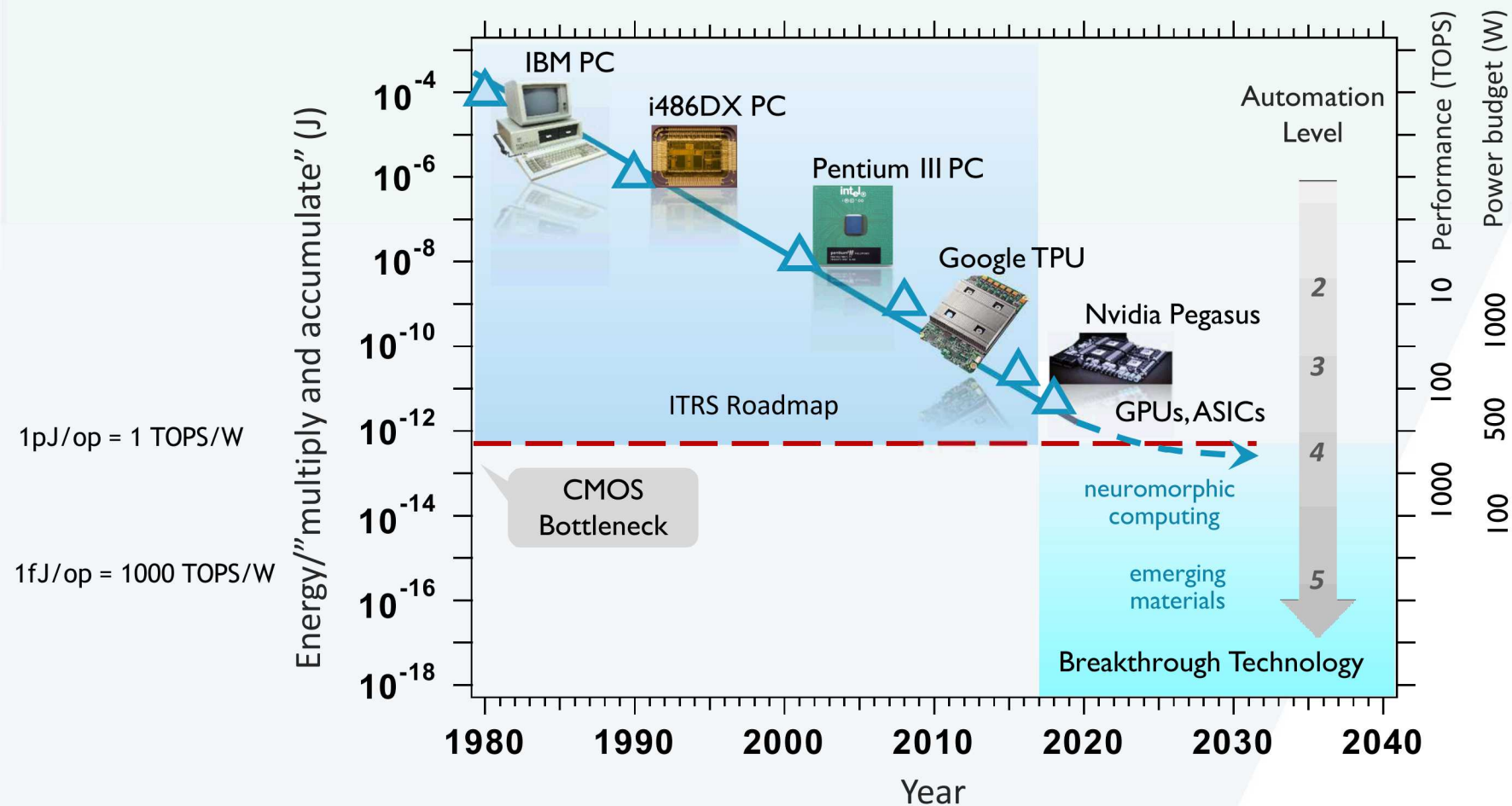
## FOCUSED RESEARCH AREAS FOR HIGHLY AUTOMATED VEHICLES

Sensors and Sensor Processing	Scene Perception and Algorithms	Navigation Hardware Accelerators
<ul style="list-style-type: none"><li>• Create disruptive optical sensing technology to reduce energy consumption by 100X</li><li>• Develop chip-scale LiDAR to reduce cost by 100X</li></ul>	<ul style="list-style-type: none"><li>• Explore sparse coding and reduced-precision to reduce computation load by 1000X</li><li>• Develop biologically inspired machine learning algorithms to reduce the number of training samples by 100X</li><li>• Develop unsupervised and self-supervised learning algorithms</li></ul>	<ul style="list-style-type: none"><li>• Develop and demonstrate hardware capable of real-time processing of tera- to petabit inputs, with energies at <math>&lt;10</math> fJ per operation (<math>&gt;100</math> TOPS/W)</li><li>• Enable algorithms for robust and reliable recognition tasks needed for perception.</li><li>• Demonstrate the value of algorithm and hardware co-design such that combined elements have greater energy and/or SWaP improvement</li></ul>



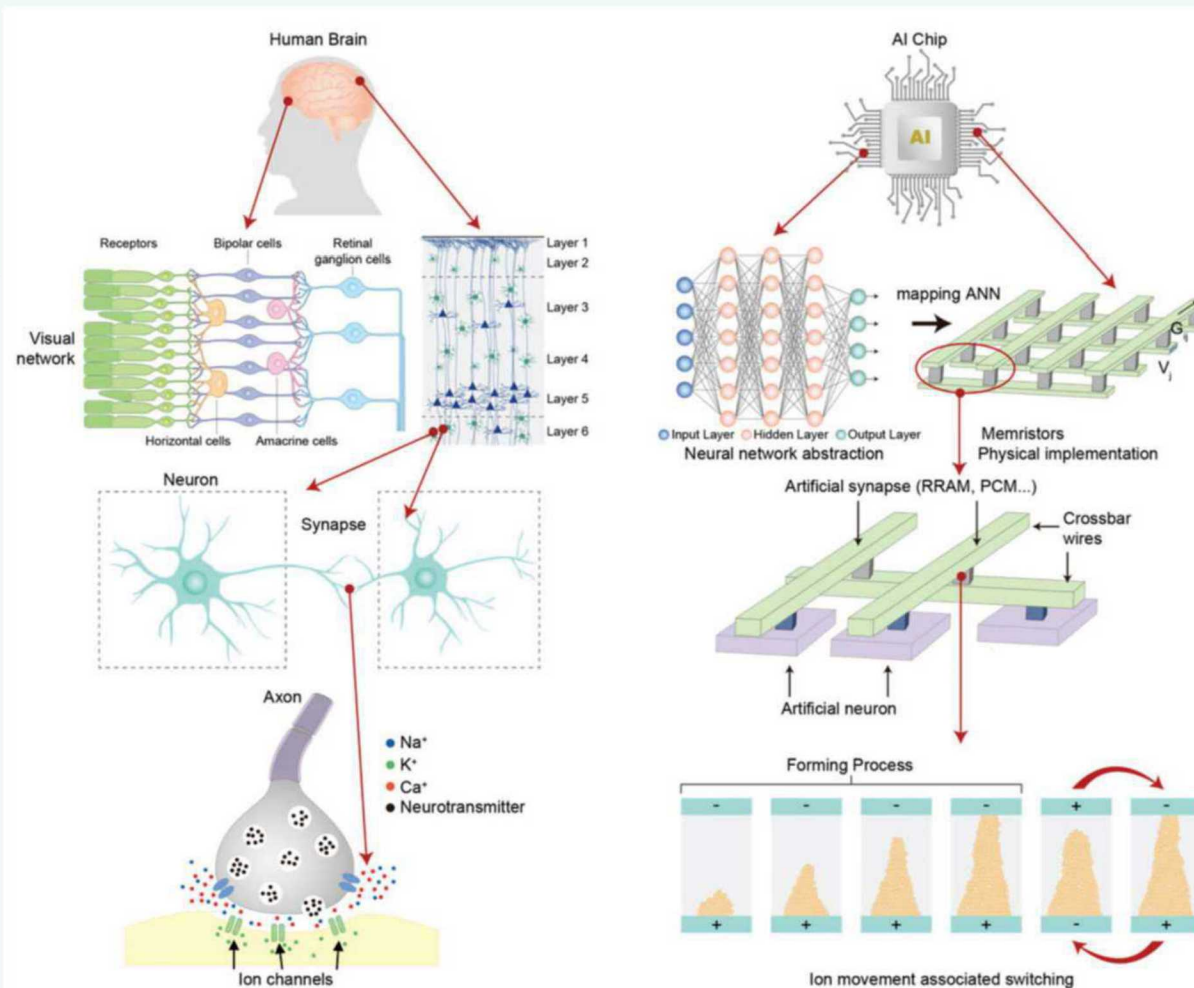


## ENERGY EFFICIENCY COMPUTING AND THE NEED OF CAV





# LEARNING FROM BRAIN FOR ULTIMATE POWER EFFICIENT COMPUTE



- Highly parallel neuron net
- Spiking - Event based computing
- Low voltage - 100mV
- Synapse - plasticity
- In memory computing with NVM (synapse)
- Matrix operation



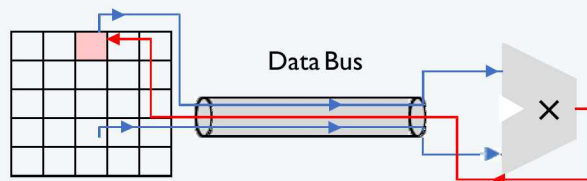
# BREAKING VON NEUMANN BOTTLENECK – UNLEASH 1000X POWER PERFORMANCE with NVM CROSSBAR

## Von Neumann Digital

Separate logic and memory structures

SRAM to store the weights

Arithmetic logic unit  
for multiplication



Uses established CMOS technology

Data bus results in latency and power

## In-memory Parallel Analog

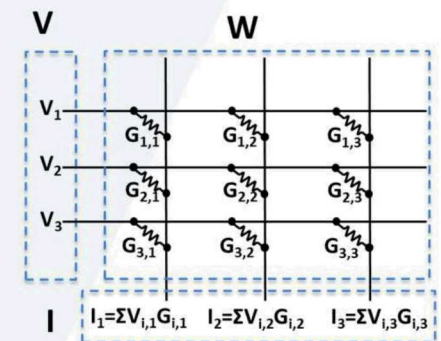
Use non-volatile memory

### Mathematical

$$V^T W = I$$

$$\begin{bmatrix} V_1 & V_2 & V_3 \end{bmatrix} \begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ W_{2,1} & W_{2,2} & W_{2,3} \\ W_{3,1} & W_{3,2} & W_{3,3} \end{bmatrix} = \begin{bmatrix} I_1 = \sum V_{i,1} W_{i,1} & I_2 = \sum V_{i,2} W_{i,2} & I_3 = \sum V_{i,3} W_{i,3} \end{bmatrix}$$

### Electrical



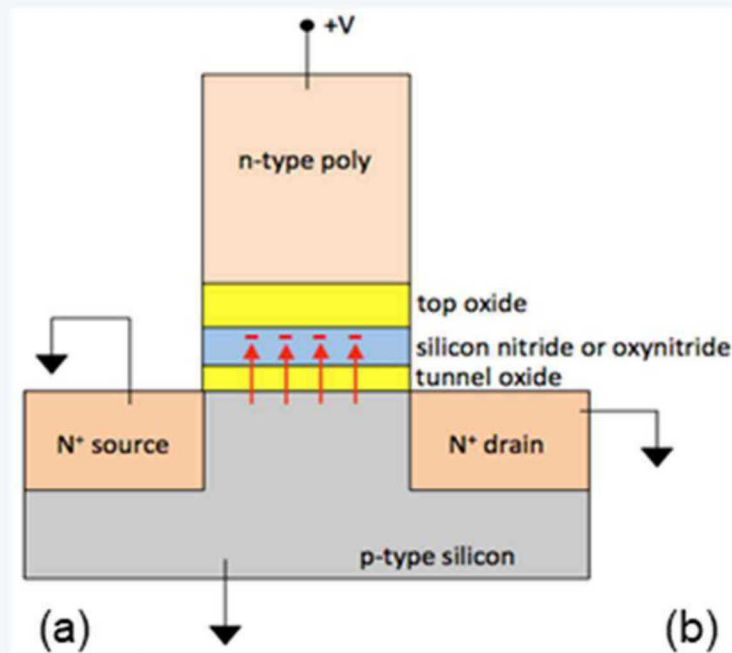
Conductance of each  
element can be changed in a  
predictable manner

Simultaneous logic and memory  
3 orders of magnitude less power

M. Marinella, *IEEE Circuits and Systems*, 8, 86-101, 2018  
Zidan, Strachan, & Lu, *Nat. Elec. I*, 22, 2018

Emerging on-chip non-volatile memory (NVM) improves energy-efficiency by performing analog multiply-accumulate inside memory and eliminate data movement

# NEAR-TERM TECHNOLOGY – 10TOPS/W for DNN accelerator



Component	Vector Matrix Multiply	Matrix Vector Multiply	Outer Product Update
Energy/Op SONOS (fJ)	13.7	13.7	68.2
Energy/Op SRAM (fJ)	2718	4630	4102
Array Latency SONOS ( $\mu$ s)	0.40	0.40	20
Array Latency SRAM ( $\mu$ s)	4	32	8

Silicon-Oxide-Nitride-Oxide-Silicon(SONOS)

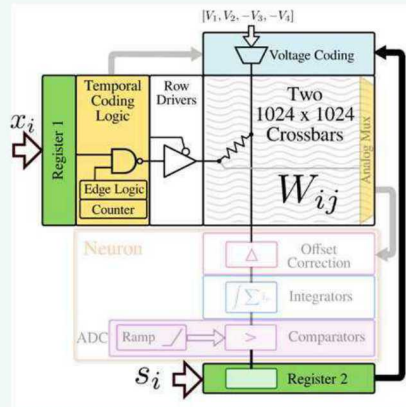
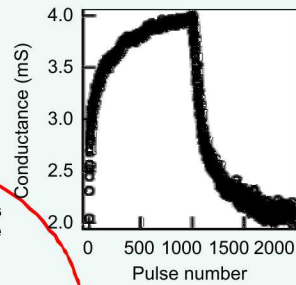
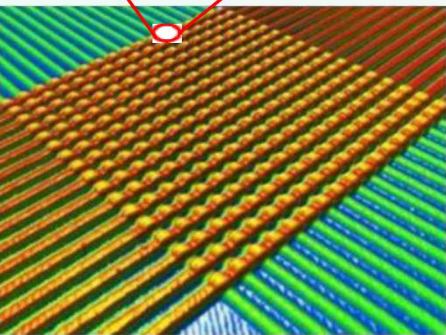
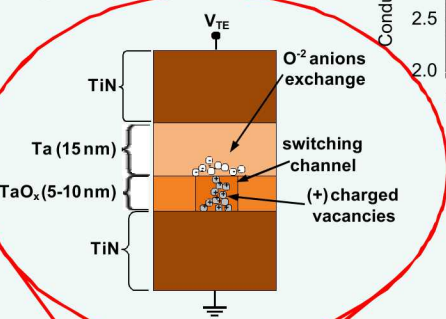
S Agarwal et al, IEEE J Exploratory Solid-State Computational Devices and Circuits, 5, 52-57, 2019.





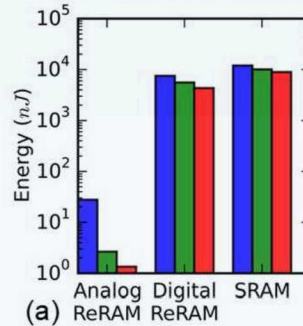
# ANALOG RERAM CROSSBAR – towards 100TOPS/W for DNN accelerator

## Memristor ReRAM

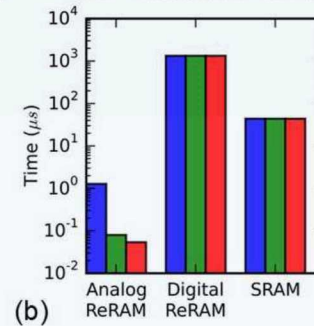


Component	Vector Matrix Multiply	Matrix Vector Multiply	Outer Product Update
Energy/Op ReRAM (fJ)	12.2	12.2	2.1
Energy/Op SRAM (fJ)	2718	4630	4102
Array Latency ReRAM (μs)	0.38	0.38	0.51
Array Latency SRAM (μs)	4	32	8

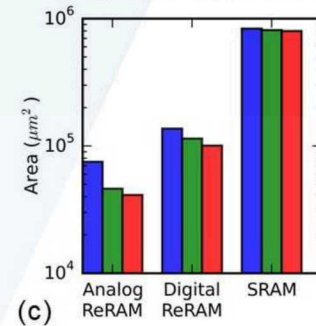
Energy  
430 – 6,900X over SRAM



Latency  
35 – 800X over SRAM



Area  
11 – 20X over SRAM



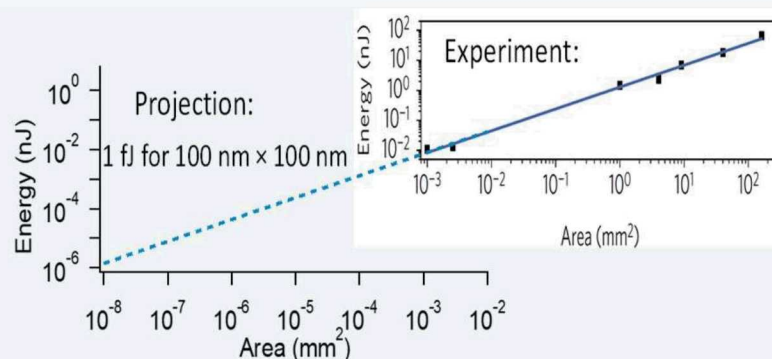
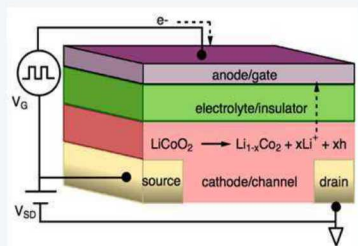
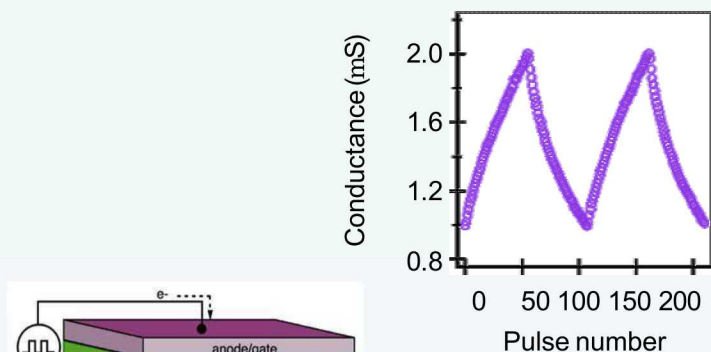
8 bit in/out 8 bit weights    4 bit in/out 8 bit weights    2 bit in/out 8 bit weights

Marinella, Agarwal, et al, *IEEE JETCAS*, 2018

Agarwal, et al, *IEEE E3S Symp*, 2017



# ION TUNABLE ELECTRONIC MATERIALS – beyond 100TOPS/W



**nature materials** LETTERS  
PUBLISHED ONLINE: 20 FEBRUARY 2017 | DOI: 10.1038/NMAT4856

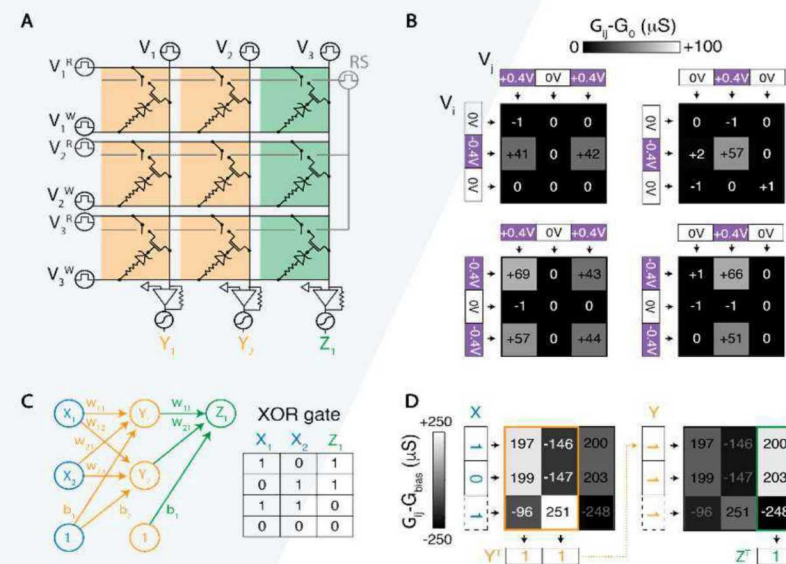
**A non-volatile organic electrochemical device as a low-voltage artificial synapse for neuromorphic computing**

Yoeeri van de Burgt<sup>1,2</sup>, Ewout Lubberman<sup>1,2</sup>, Elliot J. Fuller<sup>3</sup>, Scott T. Keene<sup>1</sup>, Grégorio C. Faria<sup>1,4</sup>, Sapan Agarwal<sup>3</sup>, Matthew J. Marinella<sup>5</sup>, A. Alec Talin<sup>1\*</sup> and Alberto Salleo<sup>1\*</sup>

**Science**

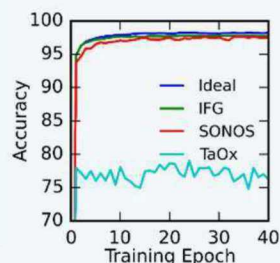
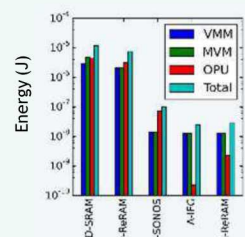
**Parallel programming of an ionic floating-gate memory array for scalable neuromorphic computing**

Elliot J. Fuller<sup>1</sup>, Scott T. Keene<sup>2\*</sup>, Armantas Melianas<sup>2\*</sup>, Zhongrui Wang<sup>2</sup>, Sapan Agarwal<sup>1</sup>, Yiyang Li<sup>1</sup>, Yaakov Tuchman<sup>2</sup>, Conrad D. James<sup>2</sup>, Matthew J. Marinella<sup>1</sup>, J. Joshua Yang<sup>2</sup>, Alberto Salleo<sup>2</sup>, A. Alec Talin<sup>1</sup>





# MULTISCALE CODESIGN FOR NEUROMORPHIC ACCELERATOR



**Energy/Performance Model**  
Model performance and energy requirements

**ROSS SIM**

**Sandia Cross-Sim:**

Translates device measurements and crossbar circuits to algorithm-level performance



**Algorithms**

**Architecture Simulation**

**Architecture**

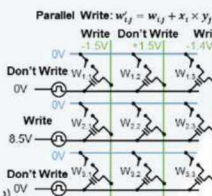
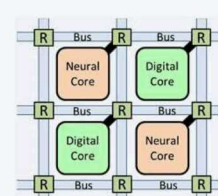
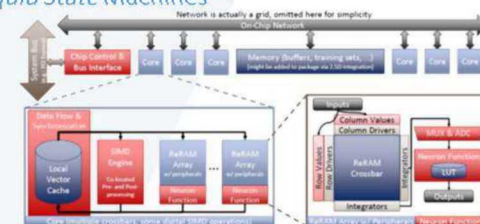
**Circuits**

**Devices**

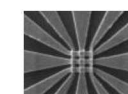
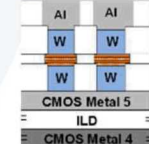
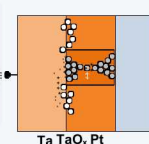
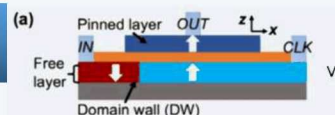
**Materials**

**Target Algorithms**

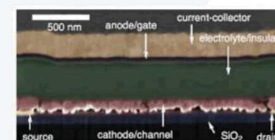
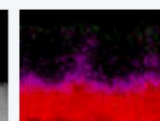
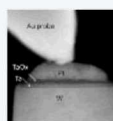
- Deep Learning
- Sparse Coding
- Liquid State Machines



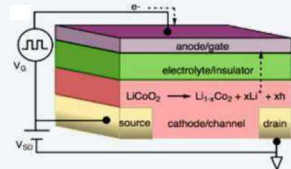
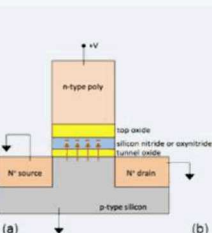
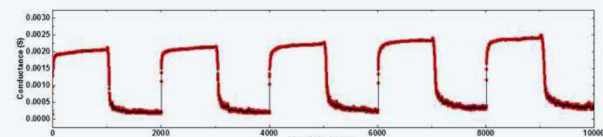
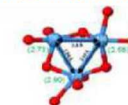
**Drift-diffusion model of transport**



**In situ Characterization**



**Ab Initio Modeling**





# COMPUTING TECHNOLOGY ROADMAPPING

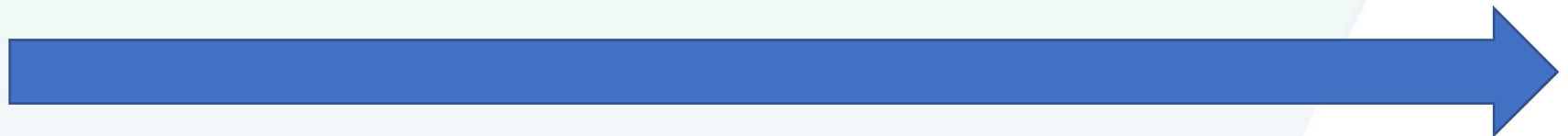
Compute perf.  
target

0.1-1 TOPS/W

1-5 TOPS/W

10 TOPS/W

100 TOPS/W



Today

2025

2035



Road-mapping

Gov't: DOE, DOT, etc.

Academic: National Labs, Research institutes

Industry: GM, Tier1, chip suppliers

- Identify the technology gaps
- Validate the target requirement for CAV
- Define the roadmap for the technology injection points





Thank you!

Q&A

Contact: Zhiyong Li  
[zli@sandia.gov](mailto:zli@sandia.gov)