



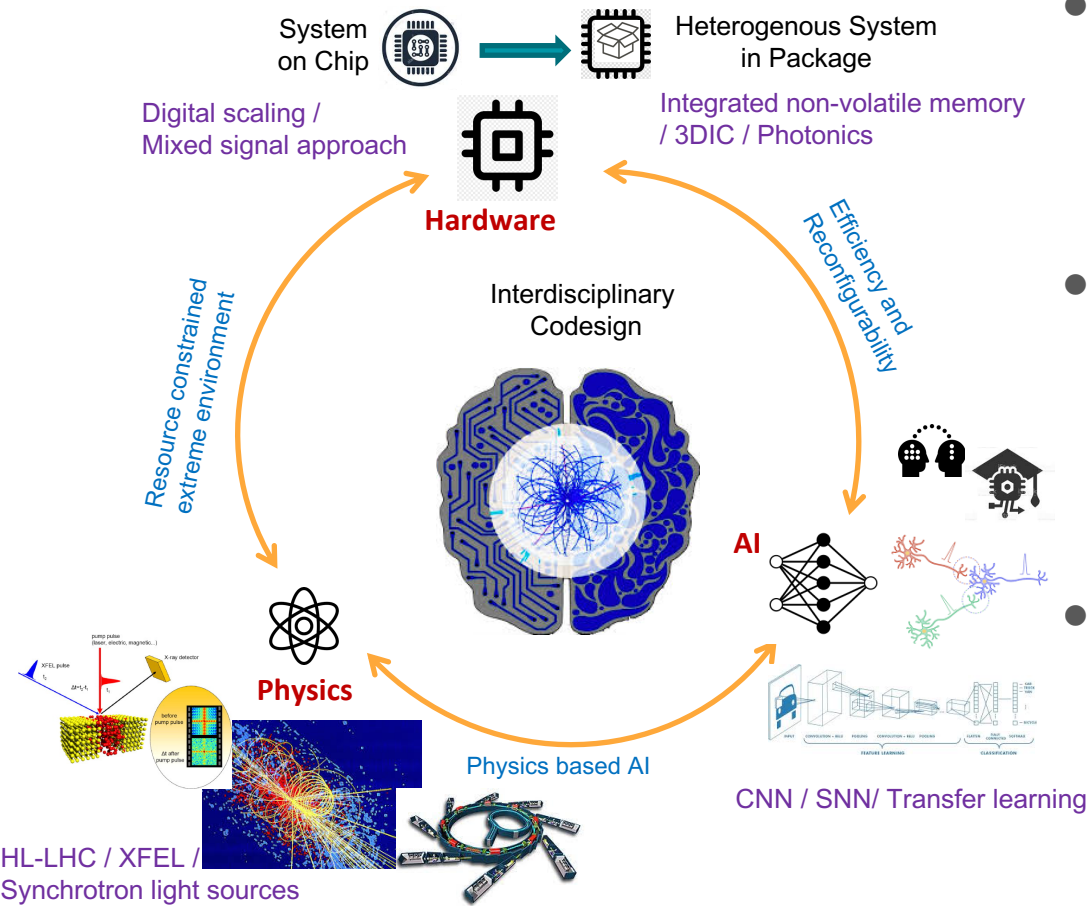
Moving intelligence onto the detector

Farah Fahim, Nhan Tran
ECFA

Why?

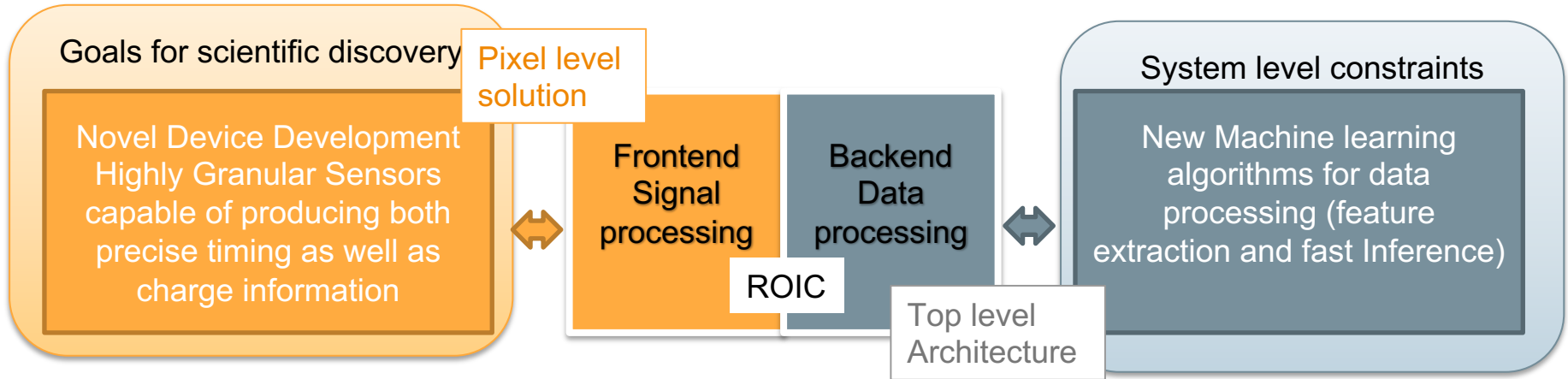
- **On/near-sensor processing valuable in reducing data rates, particularly in triggering**
 - Necessary for pixels -- pixel-triggering key capability to enable displaced physics from exotic long-lived scenarios to interesting (B)SM heavy flavor final states to dark sector physics
- **On-chip algorithm reconfigurability is powerful for dealing with changing detector and accelerator conditions**
 - Missing pixels, increasing pileup, customizable processing based on geometry
- **AI-on-sensor provides adaptable, customizable, and performant solution**
 - Configuring weights in Neural Networks provides flexible solution
 - Goal, explore traditional digital design to more speculative beyond CMOS technologies
 - Use all information, include time + space (spiking architectures)

Co-design:



- Resource constrained environment
 - High Radiation
 - Limited Power/Material budget
 - Where should this intelligence be added
- Efficiency and reconfigurability
 - Ultra-low energy per inference at extremely high rates (10's ns)
 - Reprogram both network and parameters
 - On-chip learning / inference
- Physics based Algorithms
 - Independent events
 - Depth vs. classification

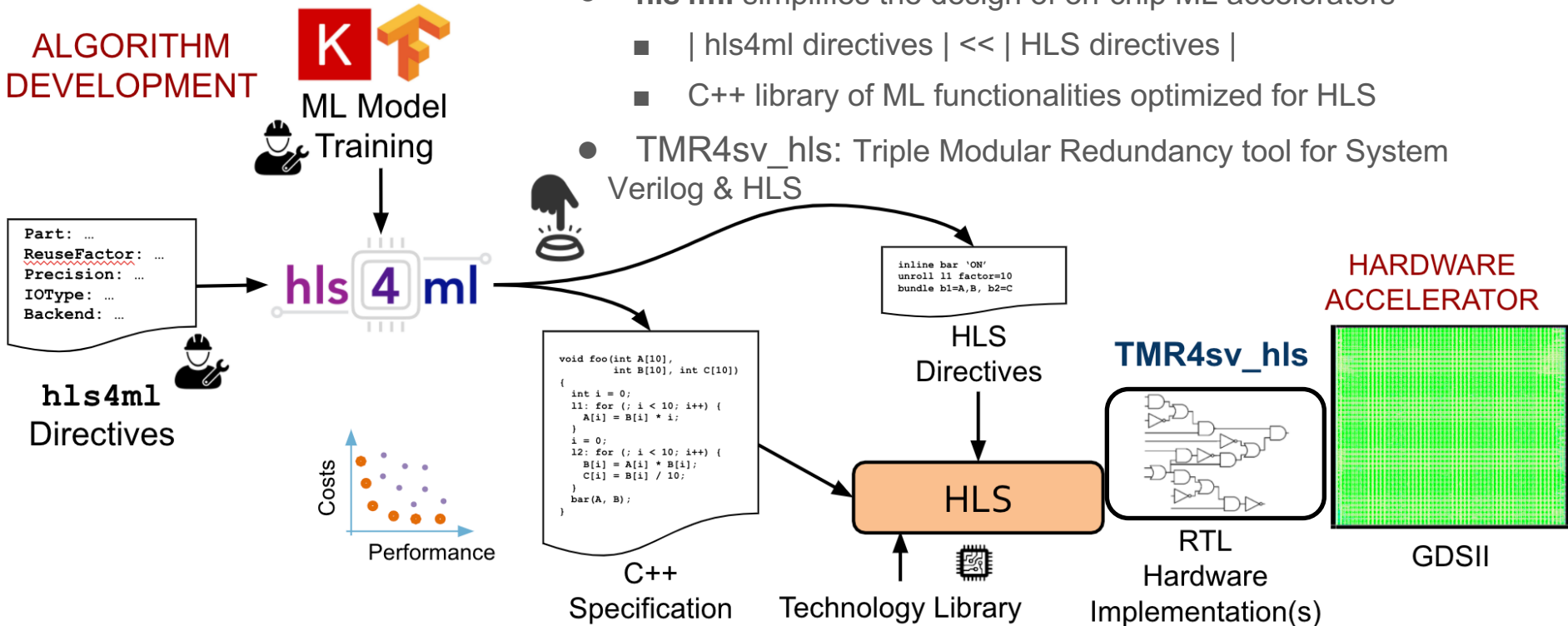
Co-design with algorithm



- Convert raw data to physics information
- Reconfigurable pixel clusters for classification dependent on detector geometries
- Create hierarchical network and enable parallel computation.

Physics Driven Hardware Co-design

- Algorithm development based on Physics data
- **hls4ml** simplifies the design of on-chip ML accelerators
 - | hls4ml directives | << | HLS directives |
 - C++ library of ML functionalities optimized for HLS
- TMR4sv_hls: Triple Modular Redundancy tool for System Verilog & HLS

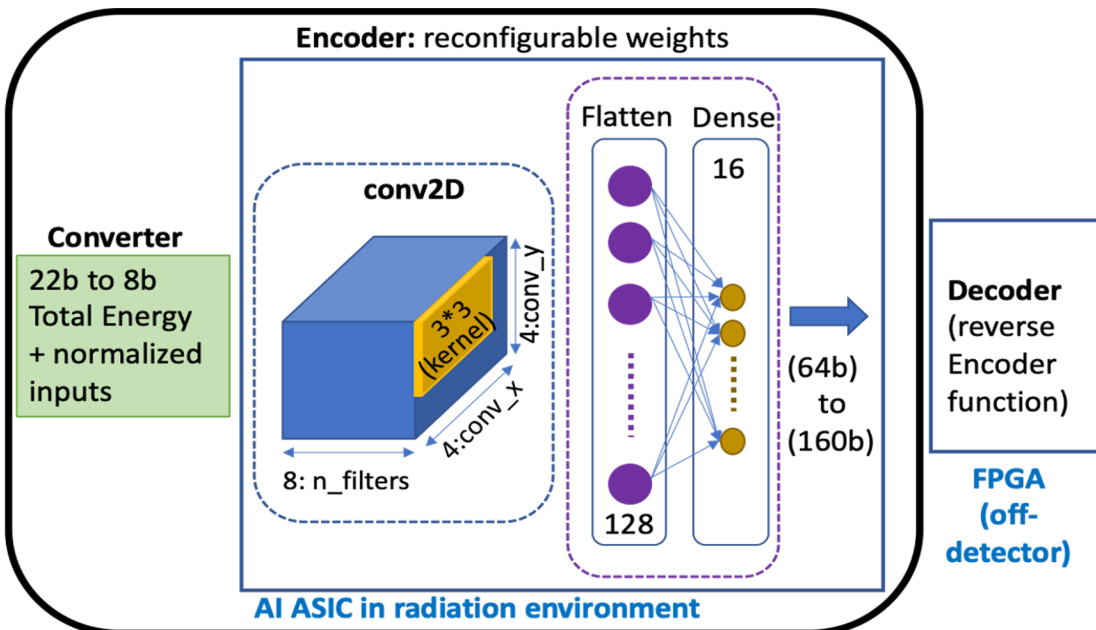
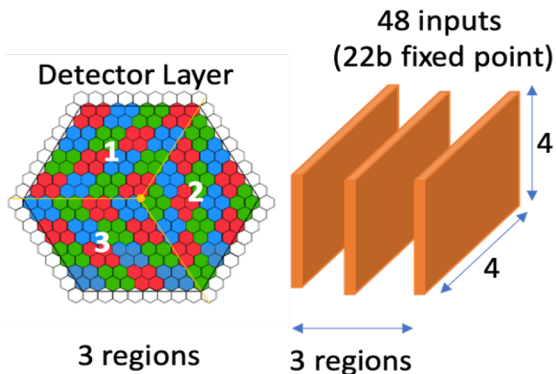


HL LHC High Granularity Calorimeter: ECON

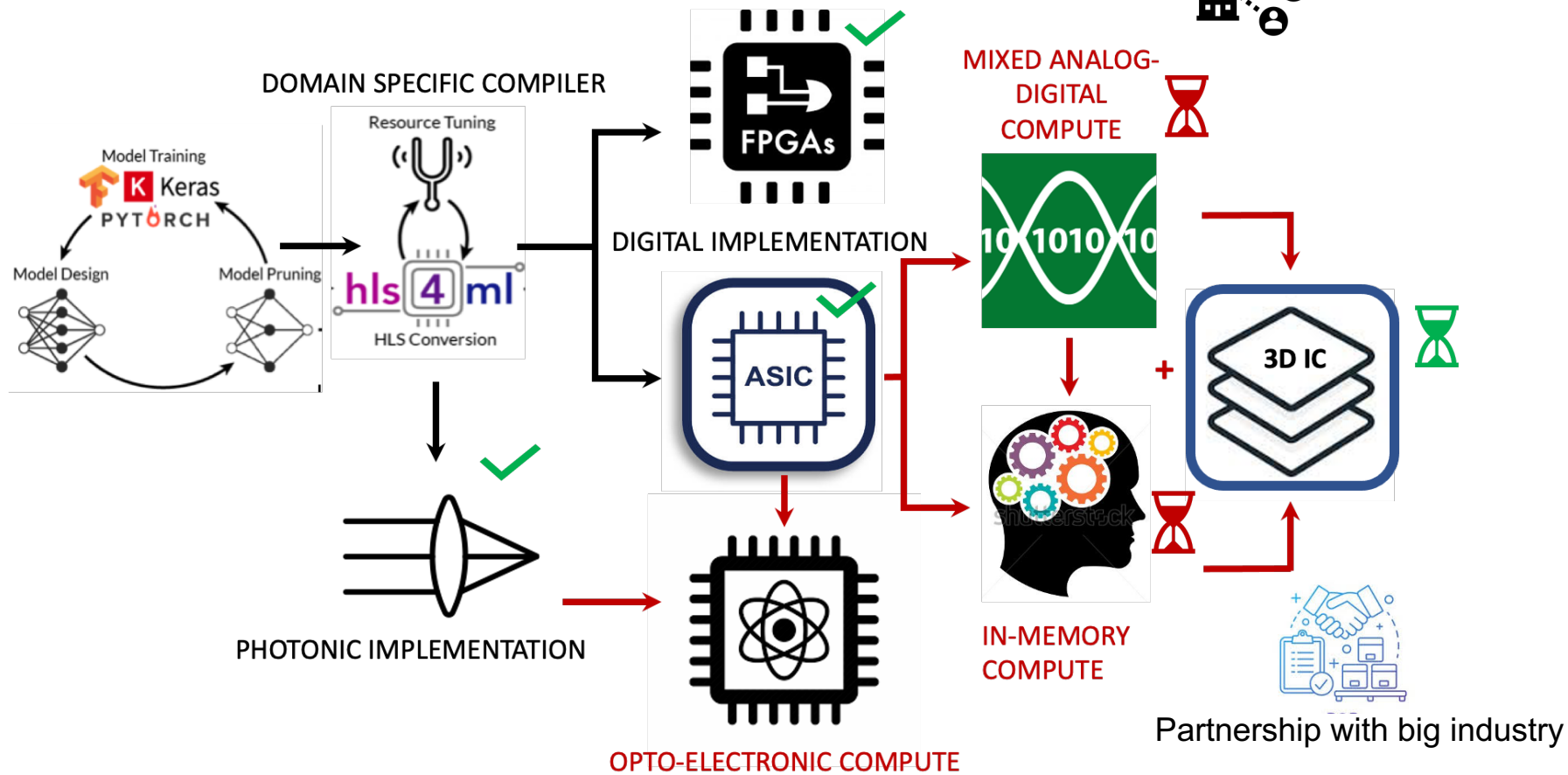
CNN: Encodes information by correlating spatial features

- **conv2D layer** – extract spatially correlated geometric features
- **Flatten layer** – Vectorizes the 2D image from the conv2D layer [8 x 4 x 4 = 128 x 1]
- **Dense layer** – aggregates the various features to provide higher order information
- **ReLU** – an activation function which introduces non-linearity by applying thresholds (part of both the conv2D and dense layers)

Metric	Simulation	Target
Power	48 mW	<100 mW
Energy / inference	1.2 nJ	N/A
Area	2.88 mm ²	<4 mm ²
Gates	780k	N/A
Latency	50 ns	<100 ns



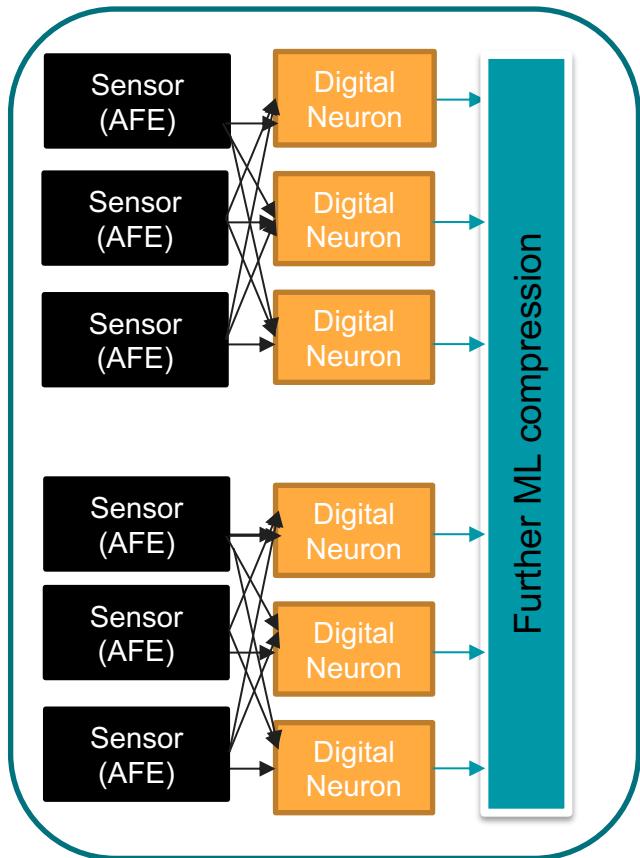
Staged approach to de-risk



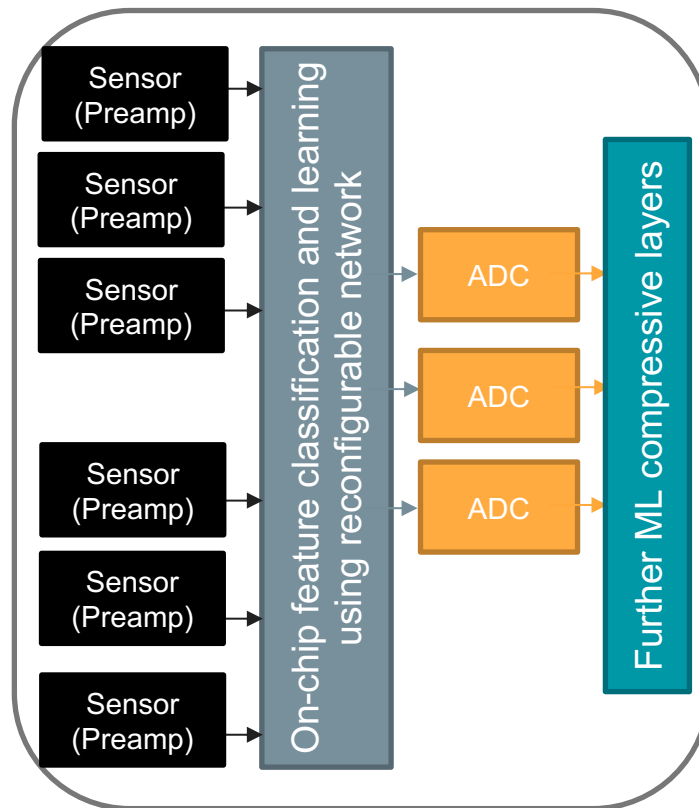
Proposed Case study: Phase III upgrade –
explore simple and novel solutions

Pixel Detector: Proposed ML implementation

Digital neuromorphic implementation

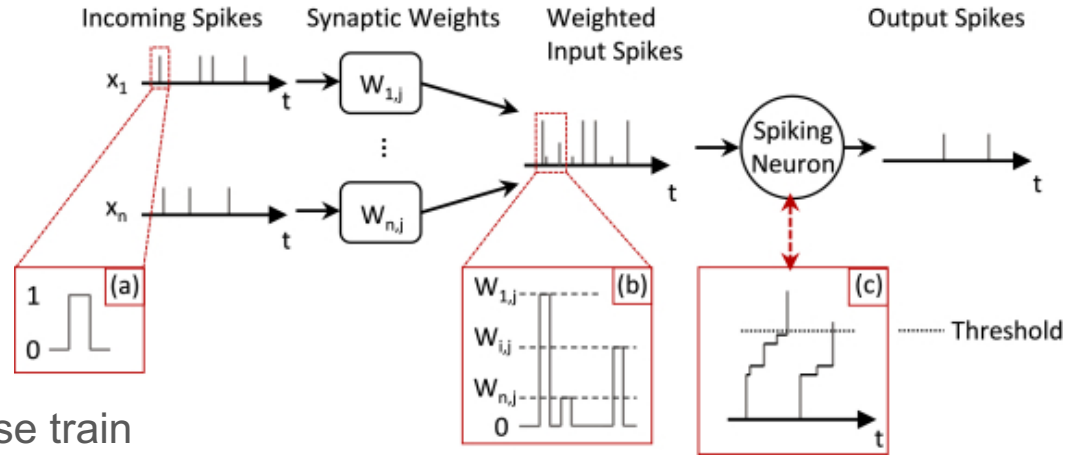


Analog – Mixed Signal implementation using floating gates or memristive cross-bar arrays



- 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?

Use digital spiking neural network

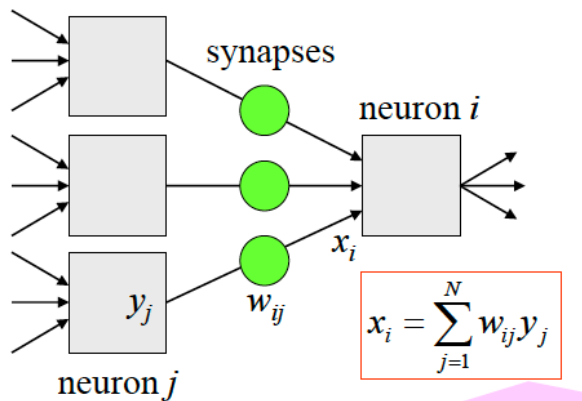


- Modify frontend to generate spikes e.g. DVS type pixel (digital pixels)
- Information is time encoded in the pulse train
- Non-uniformity correction and per pixel trimming can be handled by weight trimming
- Local neighborhood for 1st stage of classification: Compute total energy and track angle by spatio-temporal correlation
- Mature systems are based on these SNNs e.g. Loihi (14 nm), True-North, Spinnaker
- Low-power, neuromorphic approach since it runs without a clock
- HL- LHC case: since events are uncorrelated between bunches we don't need a complex network requiring historic information

Vector Matrix Multiplication

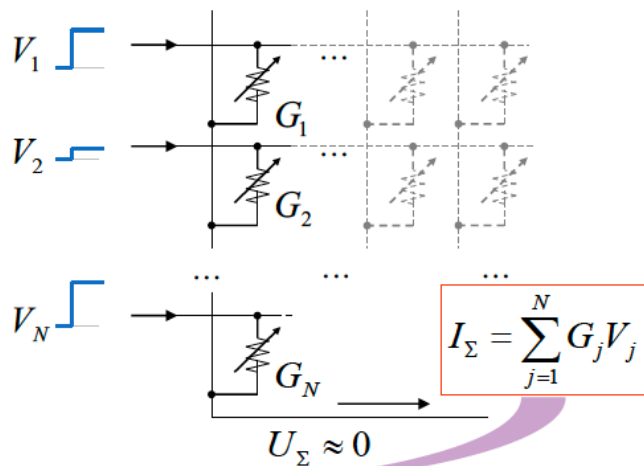
- Basic building block of neural network

Vector-by-Matrix Multiplication ...



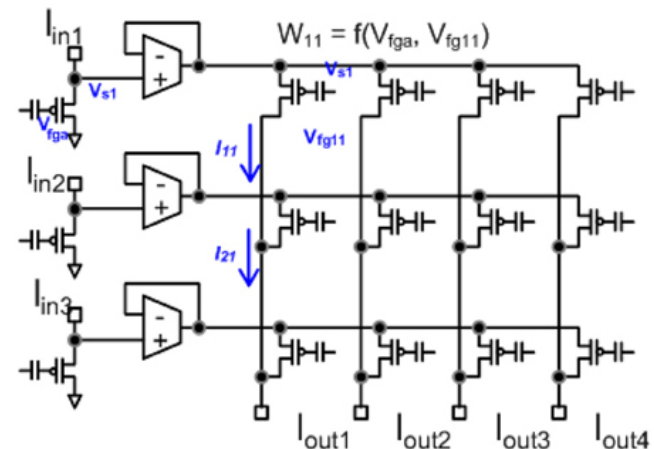
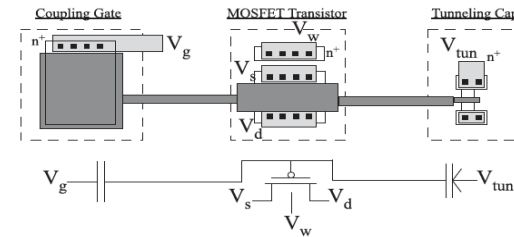
- Analog implementation
- Small footprint, programmable, large resistors

... by Analog Circuit

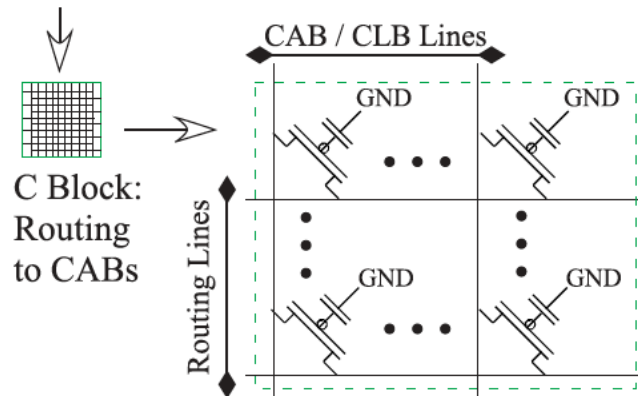


Integrate a Field programmable Analog Array (FPAA)

- Programmable Floating gate transistors for weights and switches
- Structures are available in standard CMOS process (have been demonstrated in 350 nm to 40 nm nodes)
- Radiation performance – unknown
- Uses operation transconductance amplifier (OTA) with floating gates for Vector Matrix multiplication
- Reconfigurable architecture by using a switch matrix and Manhattan routing to define interconnections



$$\begin{bmatrix} I_{out1} \\ I_{out2} \\ I_{out3} \\ I_{out4} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix} \begin{bmatrix} I_{in1} \\ I_{in2} \\ I_{in3} \end{bmatrix}$$



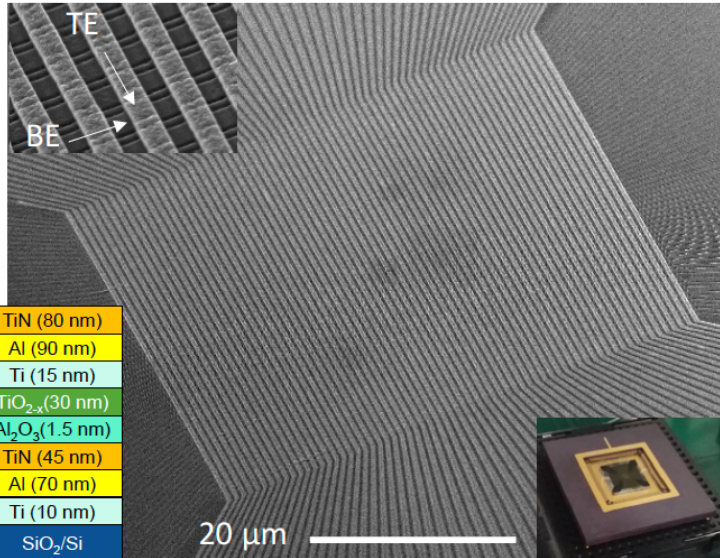
Memristive cross-bar arrays

D. Strukov, UCSB

- Use of programmable resistors (1 – 10 G Ω)
- Small footprints (< 1 μm^2)

UC Santa Barbara's Metal-Oxide Memristors

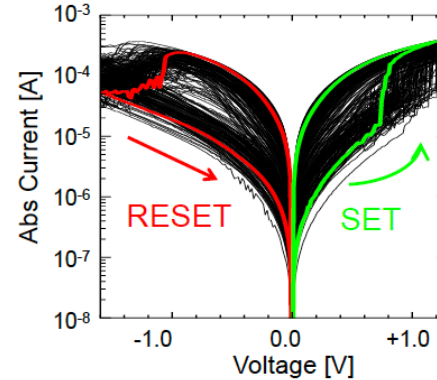
▪ 64 × 64 passive crossbar circuit



H. Kim et al. arXiv 2019

Background work: M. Prezioso et al., Nature 521, 61 2015, M. Prezioso et al. IEDM'15 p. 17.4.1, 2015, F. Merrikh Bayat et al. Nature Comm., 2018

▪ Typical I-V characteristics



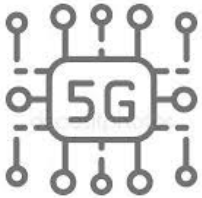
Details:

- Al₂O₃/TiO_{2-x} active bilayer by reactive sputtering
- CMOS-compatible CMP/dry etching process and TiN/Al electrodes for higher conductance
- ~250 nm wide lines
- The largest functional analog-grade passive memristor crossbar circuit supported by proper statistics

- Dense programmable memory element
- Highly energy efficient
- Foundries are also investigating non-volatile memory options (e.g. TSMC – ReRAM)
- Foundry compatible options - have highest chances of success (past experience with 3D)

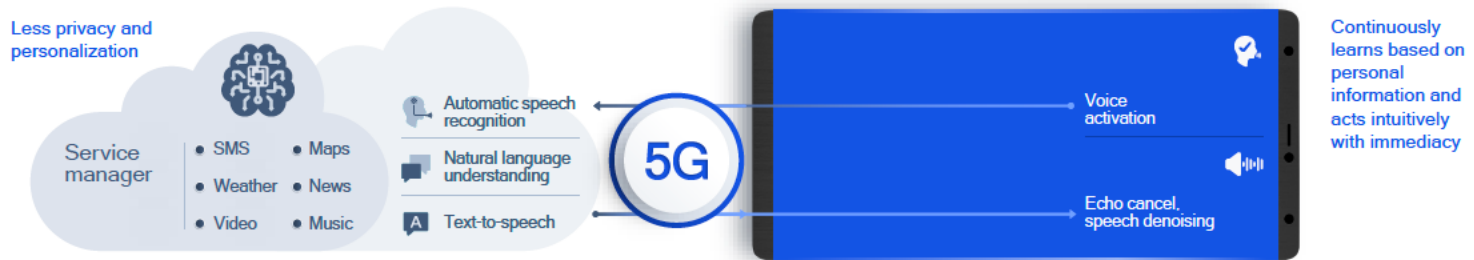
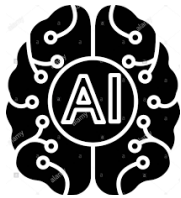
Enhance AI performance using 5G/wireless

Commercial - Qualcomm



Distributed computing enables a responsive voice UI

5G low latency allows AI tasks to be split between the device and cloud



Both ends are needed – 5G allows various implementation for appropriate tradeoffs

- Communicate between layers for more efficient data processing
- Correlations between layers can provide the best compression
- Local edge cloud can allow for low latency partial processing offload
- Use as continuous learning, additional capacity and maybe increase precision

Next steps

- Algorithm development – data sets for the various detectors
- System co-design: Programmability of parameters within power and area constraints
- Radiation studies – Floating gate capacitors / memristors
- Foundry collaborations for memristor integration (leverage growth of AI for other fields); e.g. could this become a design element in PDK like resistor/capacitor?
- Distributed learning models

Backup: Feasibility Case study: Phase II upgrade

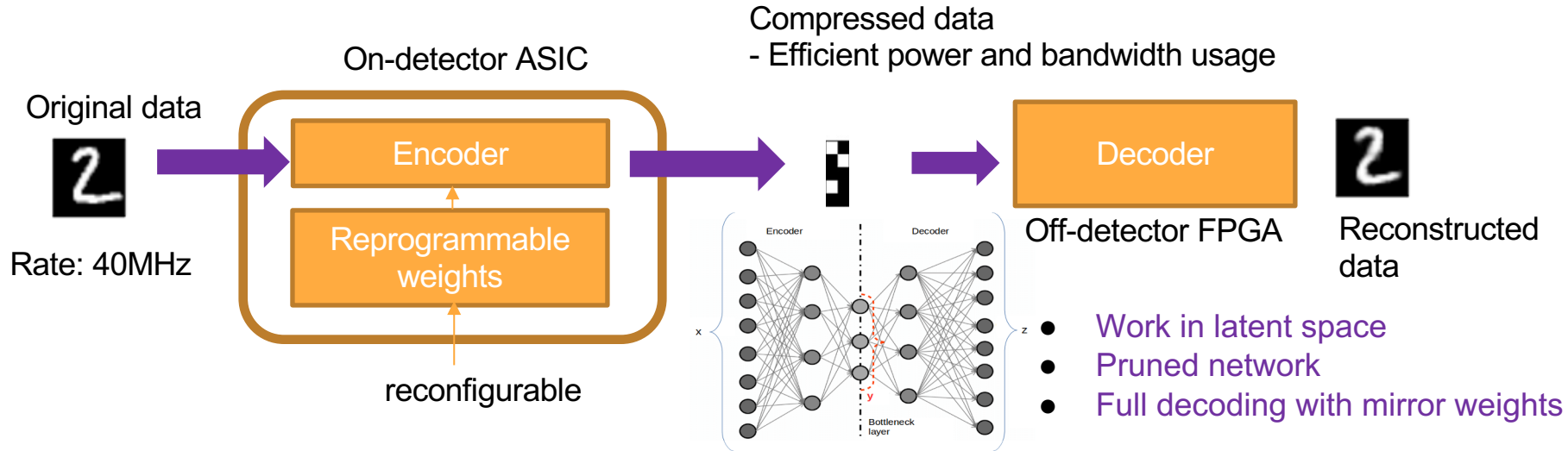
Fermilab: **Farah Fahim**, Christian Herwig, Martin Kwok, Jim Hirschauer, Nhan Tran

Northwestern University: Seda Memik, Yingyi Luo, Manuel Valentin

Columbia University: Giuseppe Di Guglielmo, Luca Carloni

Florida Tech: Danny Noonan

Deep Neural Network: Autoencoder for data-compression



- Enable edge compute : Data compression for efficient usage of power and bandwidth
- Programmable and Reconfigurable: ability to reprogram weights to adjust for detector conditions and eventually lead to self-learning intelligent detectors
- Hardware – Software codesign : Algorithm driven architectural approach
- Optimized : Low power and Low latency
- Operating in extreme radiation environment: 200 M rad
- Autoencoder for data compression is the first use case towards a DNN based on-chip learning and inference 18

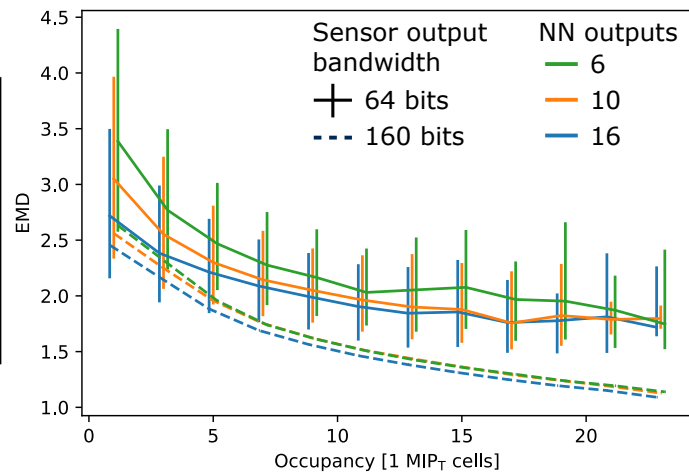
Rapid design prototyping

Neural Network architecture optimization

Lower EMD is better

Test feature	Network Architecture					Relative Power & Area		Relative Performance	
	Geometry	# filter	kernel	stride	pooling	# params	# operations	EMD Mean	EMD RMS
Reference	4x4x3	8	3x3	1	none	1.00	1.00	1.00	1.00
4x4x3 -> 8x8	8x8	8	3x3	1	none	2.73	1.76*	0.64	0.41
max pooling	8x8	8	3x3	1	2x2	0.71	0.97*	0.59	0.33
3x3 -> 5x5 kernel	8x8	8	5x5	1	2x2	0.99	2.76	0.64	0.35
pooling -> stride=2	8x8	8	3x3	2	none	0.94	0.59	0.76	0.46
8 -> 10 filters	8x8	10	3x3	2	none	1.17	0.73	0.73	0.43
8 -> 6 filters	8x8	6	3x3	2	none	0.70	0.44	0.85	0.57

* zero operations removed



Step	Type	Run Time	Iterations	Size
Model generation	D	1s	50-100	1.1k lines of C++
C Simulation	V	1s		
HLS	D	30 min	3-100	40k lines of verilog
RTL simulation	V	1 min		
Logic synthesis	D	6 hrs		
Gate-level sim	V	30 min		
Place and route	D	50 hrs		
Post-layout sim	V	1 hrs		
Post-layout parasitic sim	V	2 hrs	6	750k gates
SEE simulation	V	4 hrs		
Layout	D	20 min		
LVS and DRC	V	1 hr	1	7.6M transistors

Network optimization

Design optimization

Increasing time and complexity



Metric	Simulation	Target
Power	48 mW	<100 mW
Energy / inference	1.2 nJ	N/A
Area	2.88 mm ²	<4 mm ²
Gates	780k	N/A
Latency	50 ns	<100 ns

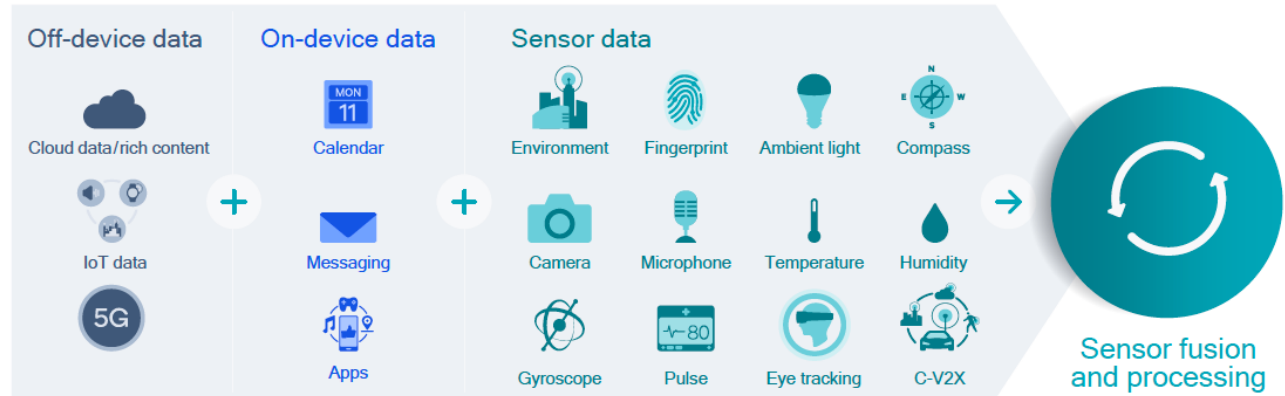
Levels of distributed AI

Data collection and processing split at 3 levels

Commercial - Qualcomm

Devices generate and possess massive amounts of data

- Sensors: Single ASICs
- Devices: Detectors / Data concentrators
- Off Devices: Processing farms



On-device AI processing of sensors and personal information conserves bandwidth while providing contextual intelligence, personalization, and privacy