# AI on chip – Algorithm to Accelerator

## Farah Fahim - Fermilab
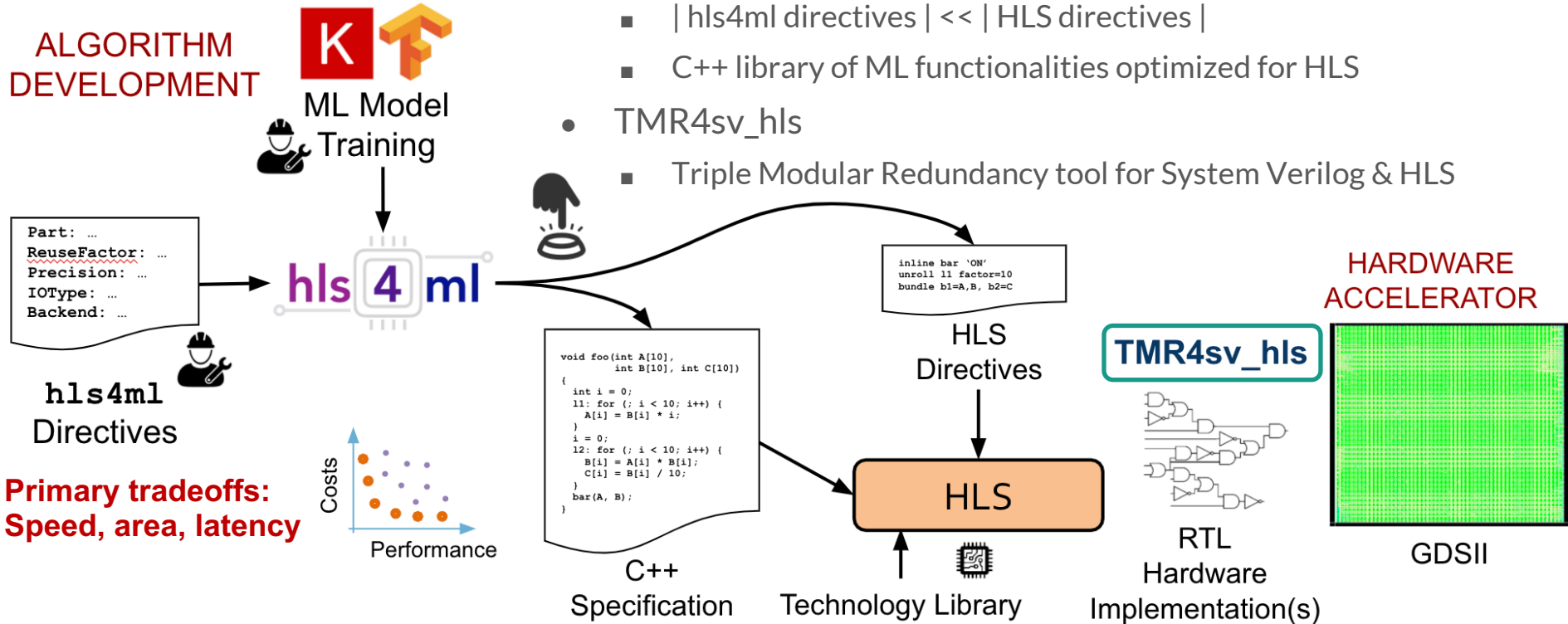
20200730

**Contributors**
Fermilab: Nhan Tran, Christian Herwig, Jim Hirschauer
Northwestern University: Seda Memik, Yingyi Luo, Manuel Valentin
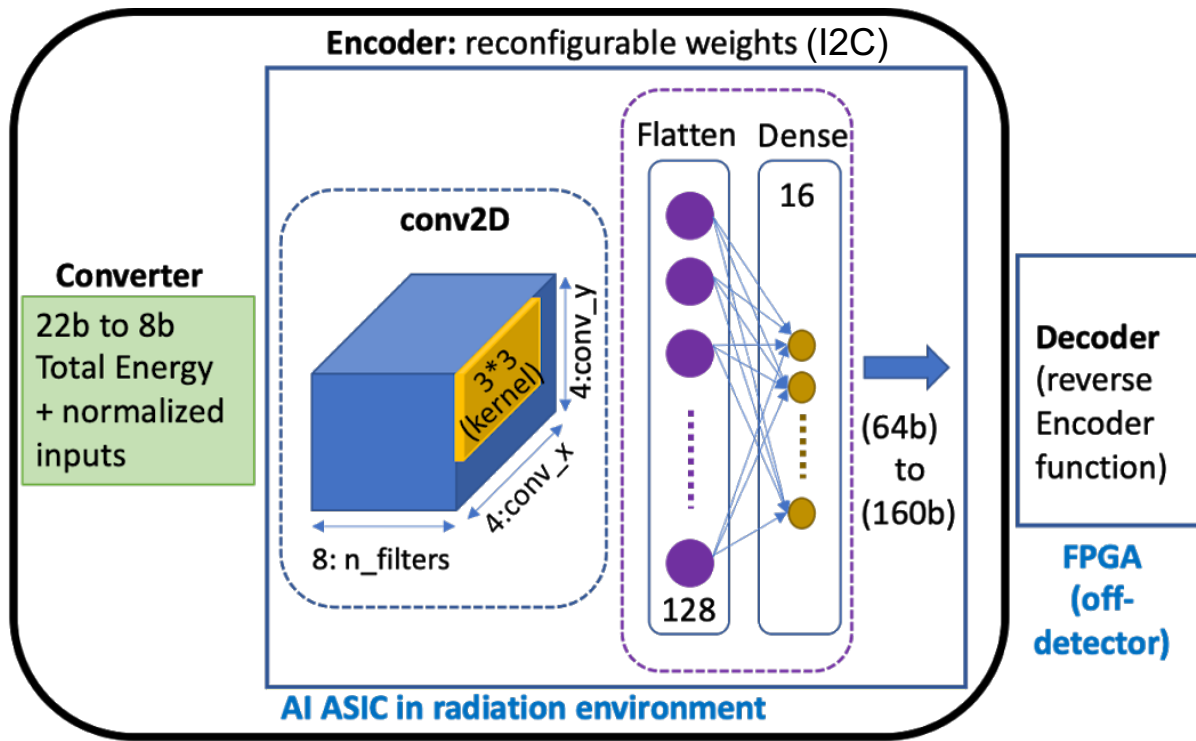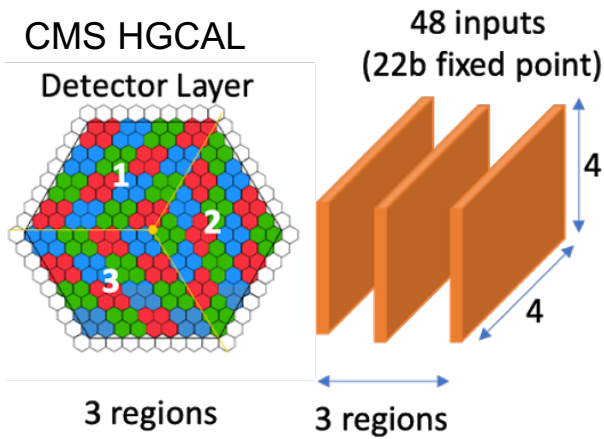Columbia University: Giuseppe DiGuglielmo

# New Paradigm for ASIC Development: 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
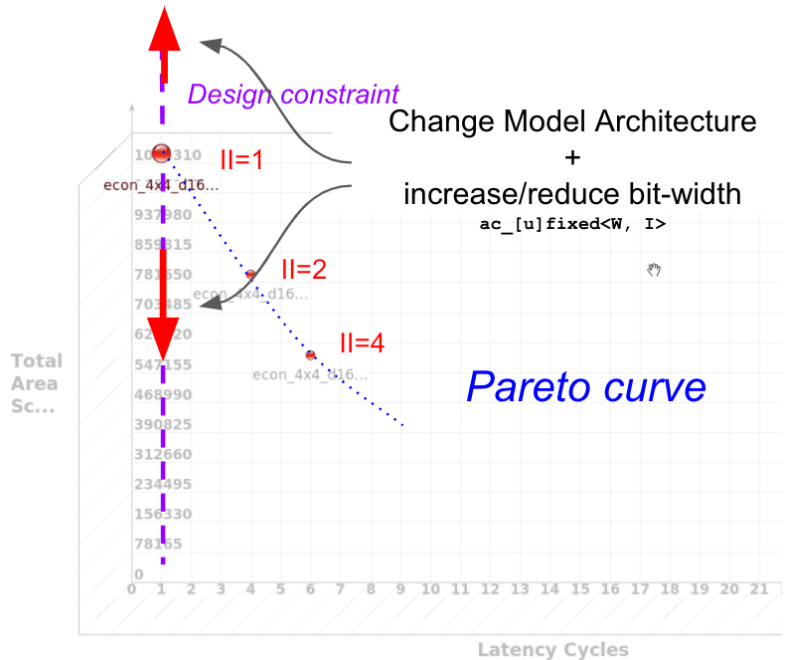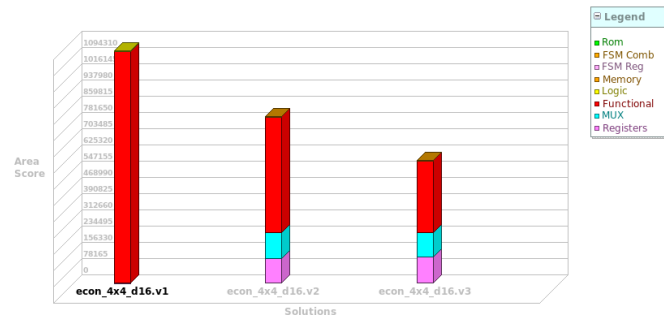  - Triple Modular Redundancy tool for System Verilog & HLS

# Autoencoder: Reconfigurable data compression

- Enable edge compute : Data compression of detector output using deep neural networks
- Programmable and Reconfigurable: ability to update weights based on **real-time feedback** (ms)
- Training **adaptable** to changing detector conditions (pileup), different geometries, lost channels, etc.

- Unsupervised learning (2.375nJ/inference, every 25ns, 15x compression, < 4mm$^2$)

# Optimization / Design Space exploration



## Report: General

| Solution | Latency Cycles | Latency Time | Throughput Cycles | Throughput Time | Slack | Total Area |
|----------|---------------|--------------|-------------------|-----------------|-------|------------|
| solution.v1 (new) | | | | | | |
| econ_4x4_d16.v1 (extract) | 1 | 25.00 | 1 | 25.00 | 13.61 | 1116589.46 |
| econ_4x4_d16.v2 (extract) | 4 | 100.00 | 2 | 50.00 | 7.54 | 802319.52 |
| econ_4x4_d16.v3 (extract) | 6 | 150.00 | 4 | 100.00 | 0.47 | 591675.33 |

Change Model Architecture
+
increase/reduce bit-width
`ac_[u]fixed<W, I>`

Design constraint

II=1

II=2

II=4

Pareto curve

Total Area Sc...

Latency Cycles

## DESIGN (D) AND VERIFICATION (V) METRICS

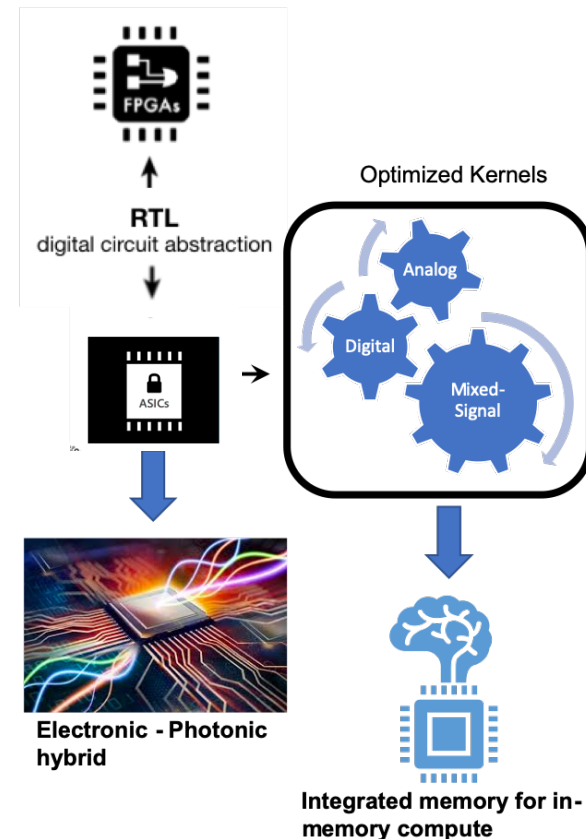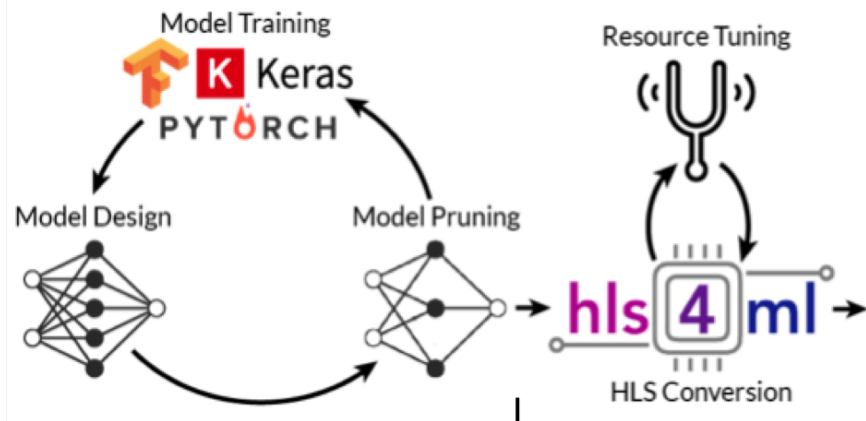| STEP | TIME | ITERATIONS | SIZE |
|------|------|------------|------|
| Model generation (D) | 0.98s | 50-100 | 1089 C++ LoC |
| C simulation (V) | 0.14s | | |
| High-level synthesis (D) | 00:30:17 | 2-3 | 39,716 Verilog LoC |
| RTL simulation (V) | 00:00:46 | | |
| Logic synthesis (D) | 06:04:19 | | 900,810 Gates |
| Gate-level simulation (V) | 00:25:19 | | |
| Place and route (D) | 71:03:53 | 6 | |
| Post-layout simulation (V) | 00:51:41 | | 1,026,387 Gates |
| Post-layout parasitic simulation (V) | 01:51:30 | | |
| Layout (D) | 00:20:00 | 1 | 12,768,389 Transistors |
| LVS & DRC (V) | 01:00:00 | | |

# Synergistic Applications
## AI on edge to AI in pixel



- Power consumption ~ 1pJ/bit - center to periphery ( ~ 5mm): routing capacitance

  - Why not do local calculations???

- **Feature extraction and data compressio**

- **Hardware driven co-design of algorithm**

# Towards heterogenous system on-chip



**OPTIMIZATION FOR POWER**
- Analog Mixed-Signal Kernels
- In-memory compute

Lowest Power Implementation of DNN
- Photonic based solution (NO Reconfigurability)

Electronic – Photonic Integrated Solution ?

Model Training
Keras
PYTORCH
Model Design
Model Pruning
hls4ml
HLS Conversion
Resource Tuning
FPGAs
RTL
digital circuit abstraction
ASICs
Optimized Kernels
Analog
Digital
Mixed-Signal

**Photonic implementation**

**Electronic - Photonic hybrid**

**Integrated memory for in-memory compute**