# PulseDL-II: A system-on-chip neural network accelerator for timing and energy extraction of nuclear detector signals

**Pengcheng Ai (Speaker)**, Zhi Deng, Yi Wang, Hui Gong, Xinchi Ran, Zijian Lang

Department of Engineering Physics, Tsinghua University

8/1/2022

# Three Elements of Design Perspectives

➢ Three independent elements:

**A.**  **Nuclear Electronics**: Readout system and signal features

**B.**  **Neural Network**: Architectural research, network training...

**C.**  **Digital Design**: NN accelerator and system-on-chip scheme

➢ Overlay of elements:

**AB.**  **Application Training**: NN algorithm research for nuclear signals, selection and optimization of network architecture

**BC.**  **Hardware Mapping**: Accelerator hardware implementation of NN

**AC.**  **System Prototype**: Hardware design in the context of readout system

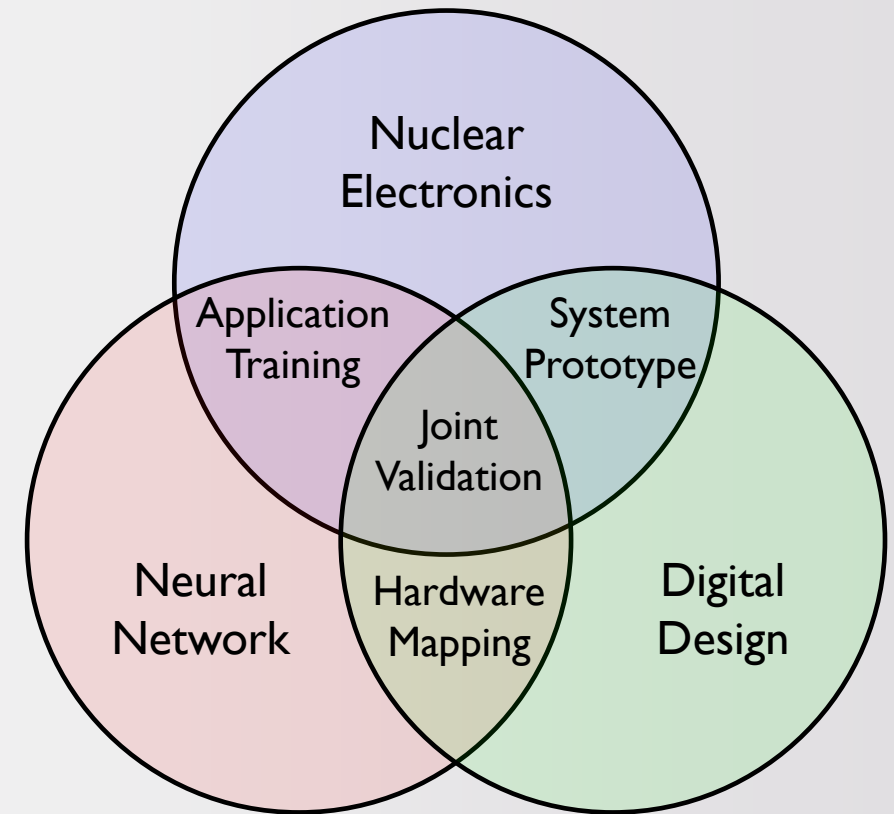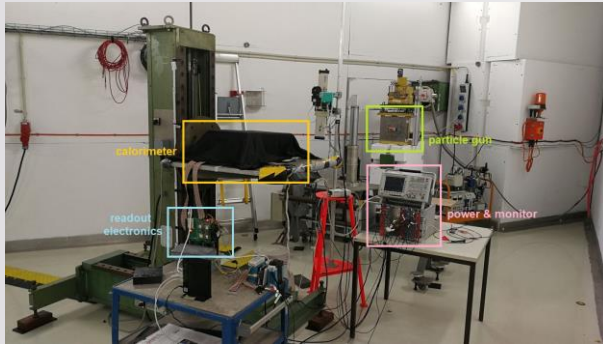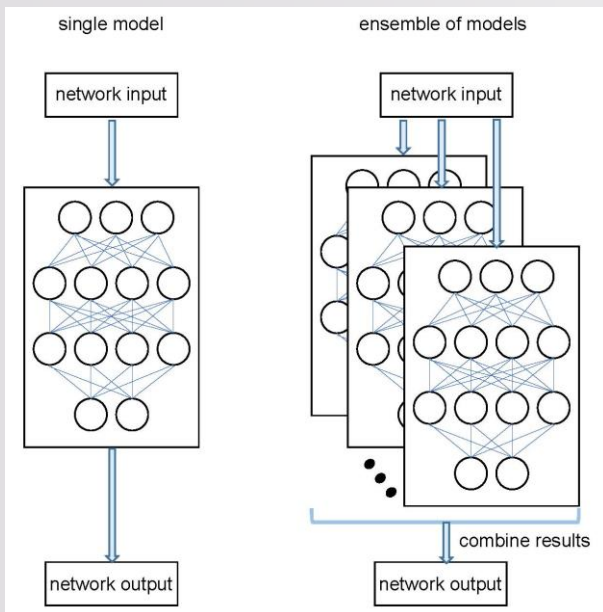**ABC**.  **Joint Validation**: Synthesis of the above three

# TABLE OF CONTENT

- Signal feature extraction with NN
- NN accelerator-based readout system
- System-on-Chip Accelerator Design
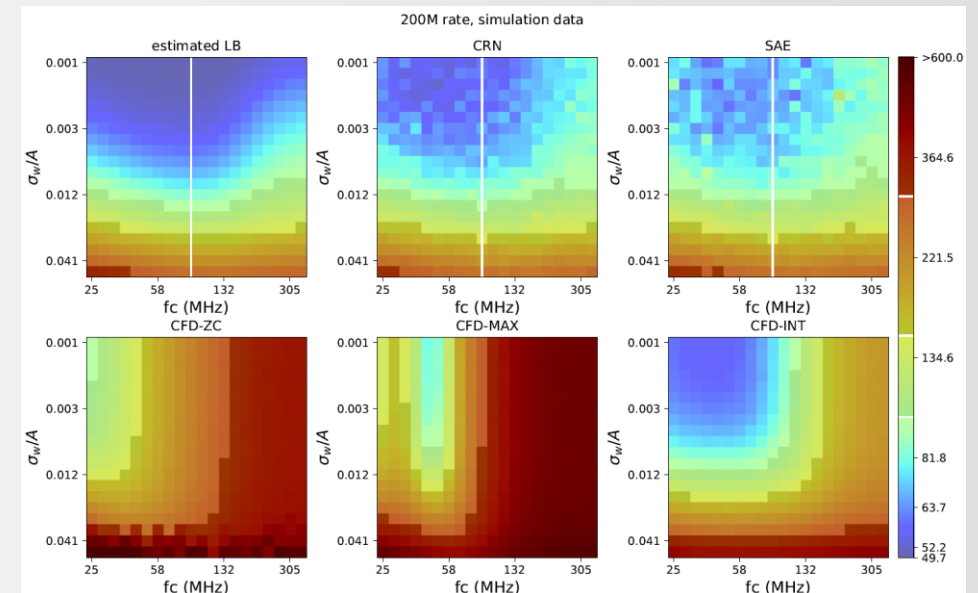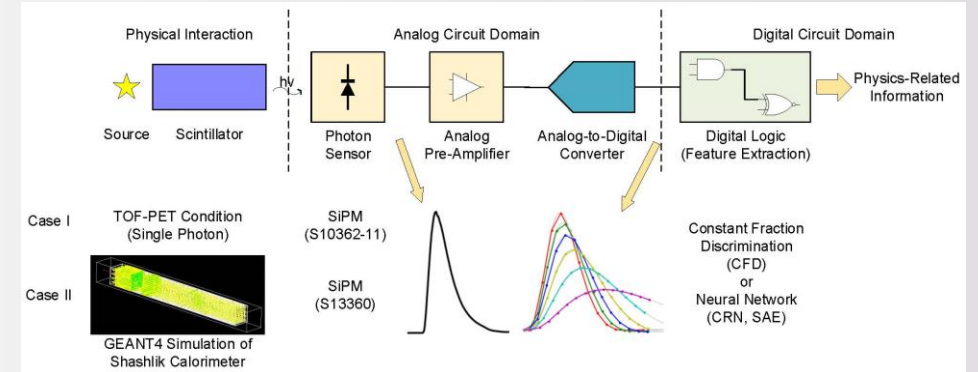- System Validation
- Summary

Estimation of heterogeneous uncertainty of nuclear detector signals with ensemble of NNs
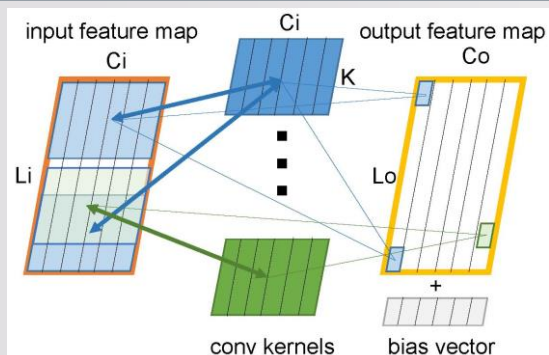
P. Ai *et al* 2022 *JINST* **17** P02032

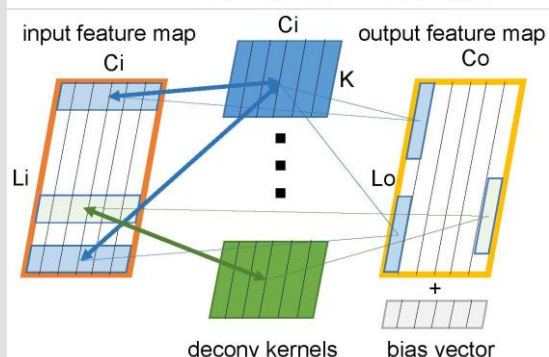Computation of the Cramer Rao lower bound of timing to find out limits for NN and traditional methods

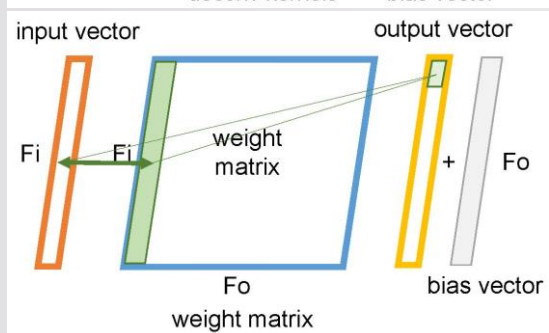P. Ai *et al* 2021 *JINST* **16** P09019
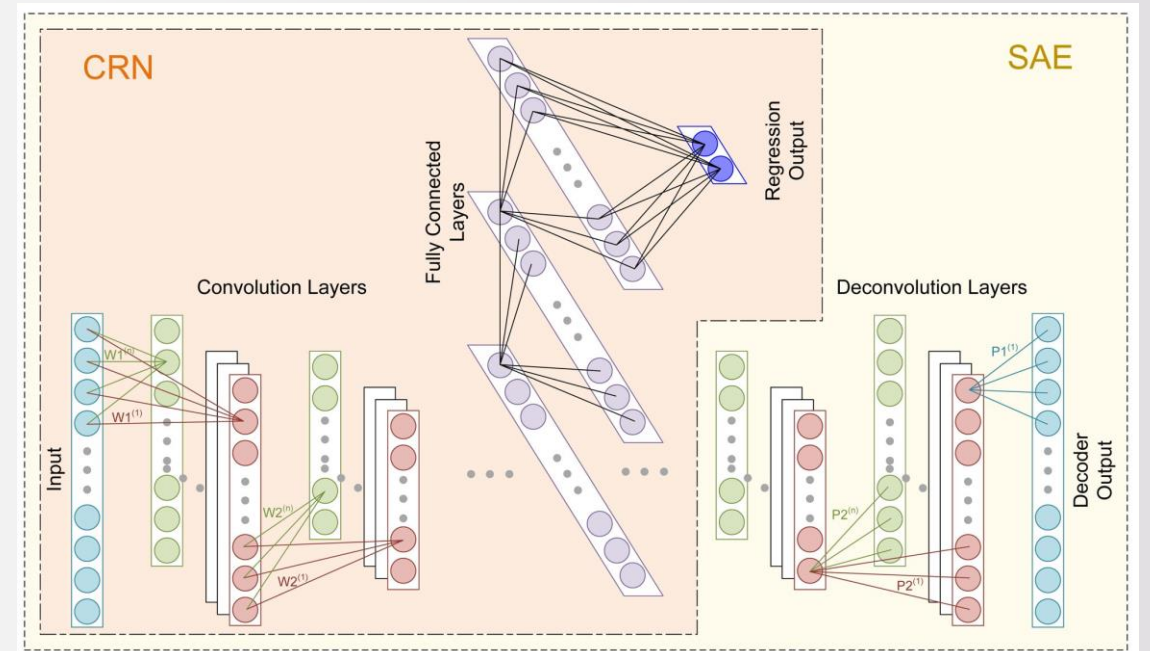
one-dimensional convolution
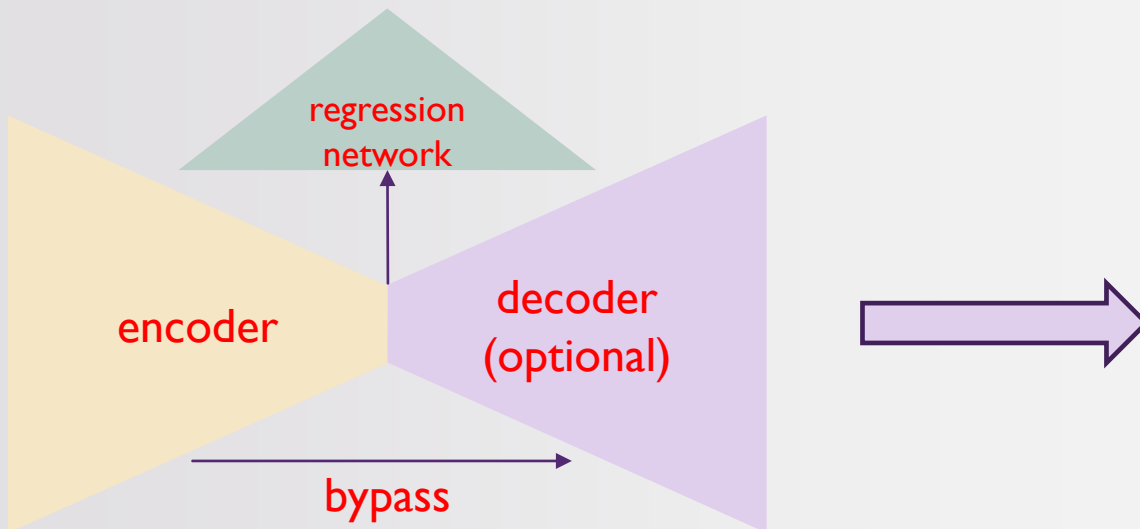
one-dimensional deconvolution

fully-connected matrix multiplication

- We choose Convolutional Neural Networks (CNN) because they succeeded in many ML tasks and facilitated parallel computing

- We select four representative building blocks:

  - 1d convolution layer

  - 1d deconvolution layer

  - fully-connected layer

  - nonlinear activation (ReLU)

- Nonlinearity is the key for **Inductive Learning** (and thus intelligent signal processing)        J.C.Ye (2022) Geometry of Deep Learning, Springer

  - Without nonlinearity, the weights in the mapping function are the same for any input sample. Once learned, they never change.        (transductive)

  - With nonlinearity, weights in the mapping function are selectively turned off/scaled by nonlinear function.        (inductive)

# Autoencoder-Based Network Architecture (AB)



- Regression network can be located at the far-end of the decoder if an accurate noiseless waveform can be obtained.
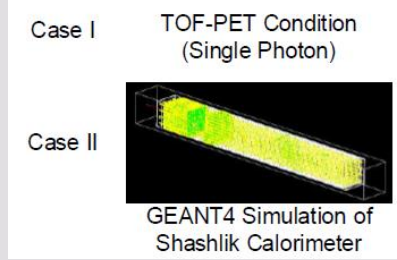
- Regression network can also be located at the bottleneck if we only have original waveform (and the decoder is optional).

Quantization-Aware Training and Validation (AB)

# Why Bringing Them to Front-End Electronics (AC)

➢ A case study of electromagnetic calorimeter (ECAL) of NICA-MPD

- 64-channel, 12-bit, 62.5M-rate ADC

- **Waveform data readout**, triggering & timing by optical fiber

- power consumption: 250 mW/channel, **water cooling system** is needed for heat dissipation

Waveform
Sampling



Fiber
Transmission

Event
Reconstruction

64-channel
ADC board

Cable
Connection


Shashlik Calorimeter

Coupling

MPPC
S13360-6025PE


SiPM

Soldering


Front-End PCB

Bias Voltage


High Voltage System

8

- A case study of electromagnetic calorimeter (ECAL) of NICA-MPD

  - 64-channel, 12-bit, 62.5M-rate ADC

  - Waveform data readout, triggering & timing by optical fiber

  - power consumption: 250 mW/channel, water cooling system is needed for heat dissipation

- Front-end upgradation with ASIC

pre-amplifier, 200M-rate ADC & NN accelerator

Shashlik Calorimeter

Coupling → SiPM

MPPC S13360-6025PE

Reduce power and bandwidth and improve performance

External Memory

Data/Instruction Transactions

Waveform Sampling

RISC Processor

Feature Output

Configuration

Feature Map

Kernel

*PulseDL* Accelerator

Fmap input

Kernel input

Temp output

Control signal

Kernel matrix

Feature map vector

Row-major order Buffer

Global control

4x4 PE array

Adder tree control

# Limitations of *PulseDL* (BC)

- The first version of the chip, although a successful practice, has the following limitations:
    - A RISC CPU outside the chip (or NN accelerator) is needed to schedule transactions
    - Dynamic quantization scheme is adopted and may bring about additional time budget
    - The adder tree structure has much space for improvement (especially the temporal adder tree)
    - Only manual configuration was done, and deep learning framework had not been supported yet


- The above limitations motivate us to develop *PulseDL-II*, the new version of the chip

- Integrate an RISC CPU into the digital design to form System-on-Chip (SoC)
    - RISC CPU: ARM Cortex-M0
    - System Bus: AHB/APB

- The *PulseDL-II* NN accelerator is mounted on the processor AHB bus as a peripheral

- Input/Output peripherals:
    - Quad/Normal SPI
    - UART (with or without internal buffer)
    - JTAG
    - GPIO

Compared to the last version:

- Adding a new topological level: Arithmetic Unit (AU)

- Broadcasting of input feature map and kernel

- Optimizing the adder tree with partial sum accumulator

- Adding function blocks for bias addition and activation

- For quantization compatible with TensorFlow or other deep learning frameworks, rescale and shift are supported

B.Jacob *et al* 2018 *CVPR* 2704  13

# Hardware-Software Codesign (BC)



**Legend:**
- Verilog/SystemVerilog
- VHDL
- Python
- C/C++
- HDF5
- HEX file

Hardware ← | → Software

**Testbench (Verilator)** ← **NN Pulse Processor**

**Quantized Simulation Data** → **Input/Output Transactions**

**Model Quantization (TensorFlow-Compatible)** ← **TensorFlow Neural Network Model**

**Model Parameters**

**Testbench (VUnit)** ← **Dual-Port AHB RAM**

**Hexadecimal Data Files**

**C Headers** **CPU Software**

**Cortex-M0 SoC** **Hexadecimal Program File** **ARM GCC Compiler**

**SoC Testbench (VUnit)**

**HDL** to design hardware components

**Neural Network Related**

**Verilator & VUnit** for simulation

**Python** to develop NN in TensorFlow framework

**Accelerator System Related**

**HDF5** to save design database

**C** to design ARM embedded software

14

# Embedded Software with Weight-Stationary Mapping (BC)

- The designed hardware allows different mapping rules
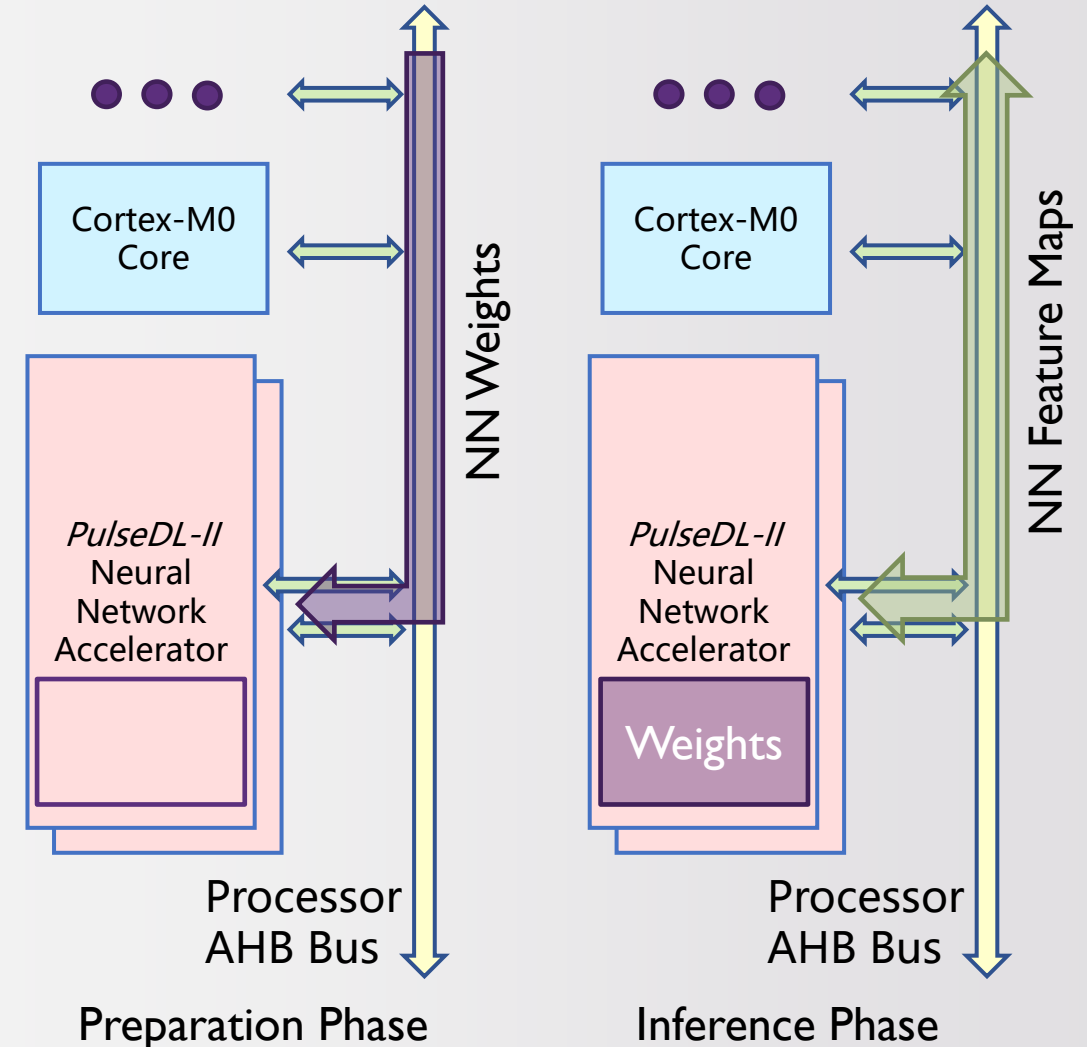
- For NNs with small/medium size, a **weight-stationary** mapping scheme can be adopted

  - Weights are stored into PEs before samples come in (Preparation Phase)

  - Only input data, output data and intermediate feature maps are transferred during inference (Inference Phase)

- The embedded software enables following features:

  - **Layer-wise inference pipelining**: weights for different layers are mapped to different groups of PEs, and they can operate simultaneously

  - **Event-level parallelism**: Each event is assigned a unique token, which will be passed in company with feature maps along the pipeline
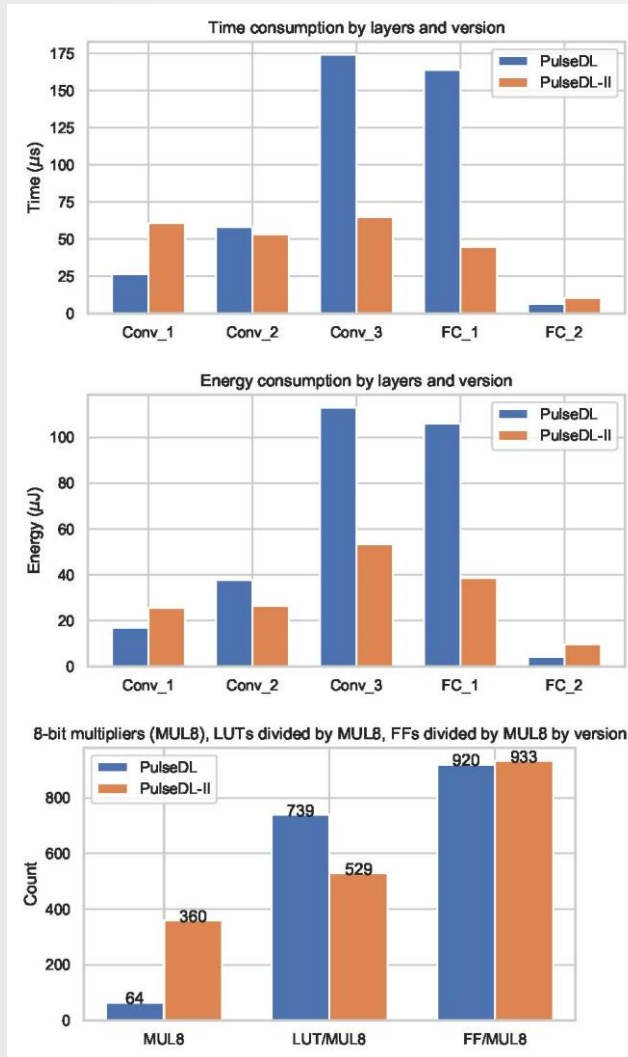


Preparation Phase

Inference Phase

15

- Evaluation settings
  - Xilinx ZCU104 Evaluation Board
  - 100 MHz working frequency
  - post-synthesis

(*PulseDL-II* NN accelerator is isolated for fair comparison with *PulseDL*):



Performance
(time consumption)
1.83x less

Power
(energy consumption)
1.81x less

Area
(resource utilization)
comparable or less



inputs (12-bit ADC)

conv layer 1
MAC #: 512 | Param #: 80

conv layer 2
MAC #: 8192 | Param #: 2080

conv layer 3
MAC #: 16384 | Param #: 8256

fc layer 1
MAC #: 8192 | Param #: 8256

fc layer 2 (outputs)
MAC #: 128 | Param #: 130

**Total**
MAC #: 33.4K | Param #: 18.8K

Used NN workload

16

Host Computer
1. FPGA Firmware, Integrated Logic Analyzer; 2. ARM MCU Program; 3. Feature Output

SDP-K1
ARM MCU

AD9106-ARDZ-EBZ
DDS Signal Generator
156.25M, 12b, 2Vpp

Ch1
Ch2
Ch3
Ch4

ADS4225 Card
125M, 12b, 2Vpp

ADS4225 Card
125M, 12b, 2Vpp

MZU07A-EV
FPGA Dev Board

ILA 1

ILA 0

VIO 0

ILA 2

Self Trigger

Trigger Logic

Monitor Adapter

BRAM

Check Point

FEP Card Interface

Ring Buffer

Time Stamp FIFO

Data FIFO

Event DAQ

Neural Network Adapter

BRAM

From ADC FEP Cards

FEP Card Wrapper

To Accelerator System-on-Chip

ILA: Integrated Logic Analyzer

VIO: Virtual Input/Output

18

**time**  **energy**

- Test waveform:



traditional methods (CFD, integral)

$$s(t) = K\left(\frac{t - t_0}{\tau}\right) e^{-\frac{t-t_0}{\tau}} u(t - t_0)$$
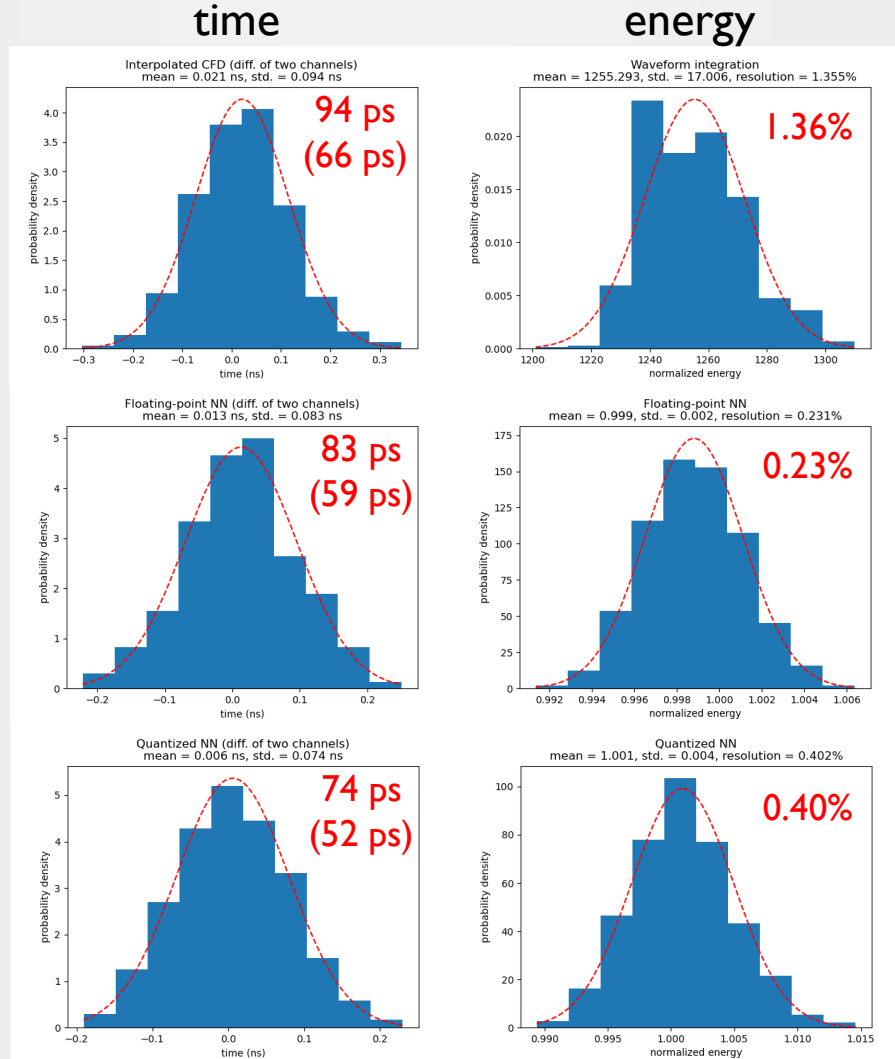
$$\tau = 40\ ns, \qquad K = K_1 K_2$$

$$SNR = 20 log_{10}\left(\frac{K_1}{\sigma_{base}}\right) = 47.4\ dB$$

$$K_2 \sim U(0.5, 2.0)$$

floating-point NNs

quantized NNs

- Sample **32** points per event

- Dual-channel synchronous waveform input



94 ps (66 ps)

1.36%

83 ps (59 ps)

0.23%

74 ps (52 ps)

0.40%

- Runtime statistics:
  - Zynq UltraScale+

| Resources (area): | |
|---|---|
| LUT | 2825 + 89540 |
| FF | 517  + 75028 |
| BRAM | 8.0   + 48.0 |
| URAM | 8     + 0 |
| **Power:** | |
| Dynamic | (0.371 + 0.541) W |
| Static | 0.594 W |
| **Performance @ 100 MHz:** | |
| Internal inf. | 113.8 us |
| Throughput | 8.3k events/sec |

# Summary

- The ability and potential of NNs in signal feature extraction are investigated

- Application-specific NN architectures are designed

- NN accelerator-based front-end electronics is prototyped

- System-on-Chip digital system with NN accelerator is developed

- System Validation on FPGA platform is done

What's next:

- Evaluate the whole system in real-world nuclear detector dataflows

- Design optimization, ASIC layout, tape-out with advanced technology

# THANK YOU!

## ANY QUESTIONS?