# Versal ACAP Processing for HL-LHC Calorimeters Signal Reconstruction

**27th Conference on Computing in High Energy and Nuclear Physics (CHEP)**
**24th October, 2024**

*Francisco Hervas*, *Alberto Valero, Luca Fiorini, Hector Gutierrez*

# Table of contents

# 1 Introduction

## *LHC calorimeter read-out*

- In the **LHC**, **Bunch Crossings** (BC) happen at **40 MHz** (25 ns)
- The **processing** happens after the **Level-1 Trigger**, at **100 kHz** (10 us)
- Signals are processed online using the **Optimal Filtering** (OF) algorithm
  - The processing is made using **Digital Signal Processors** (DSPs)
  - Therefore, it is **sequential**
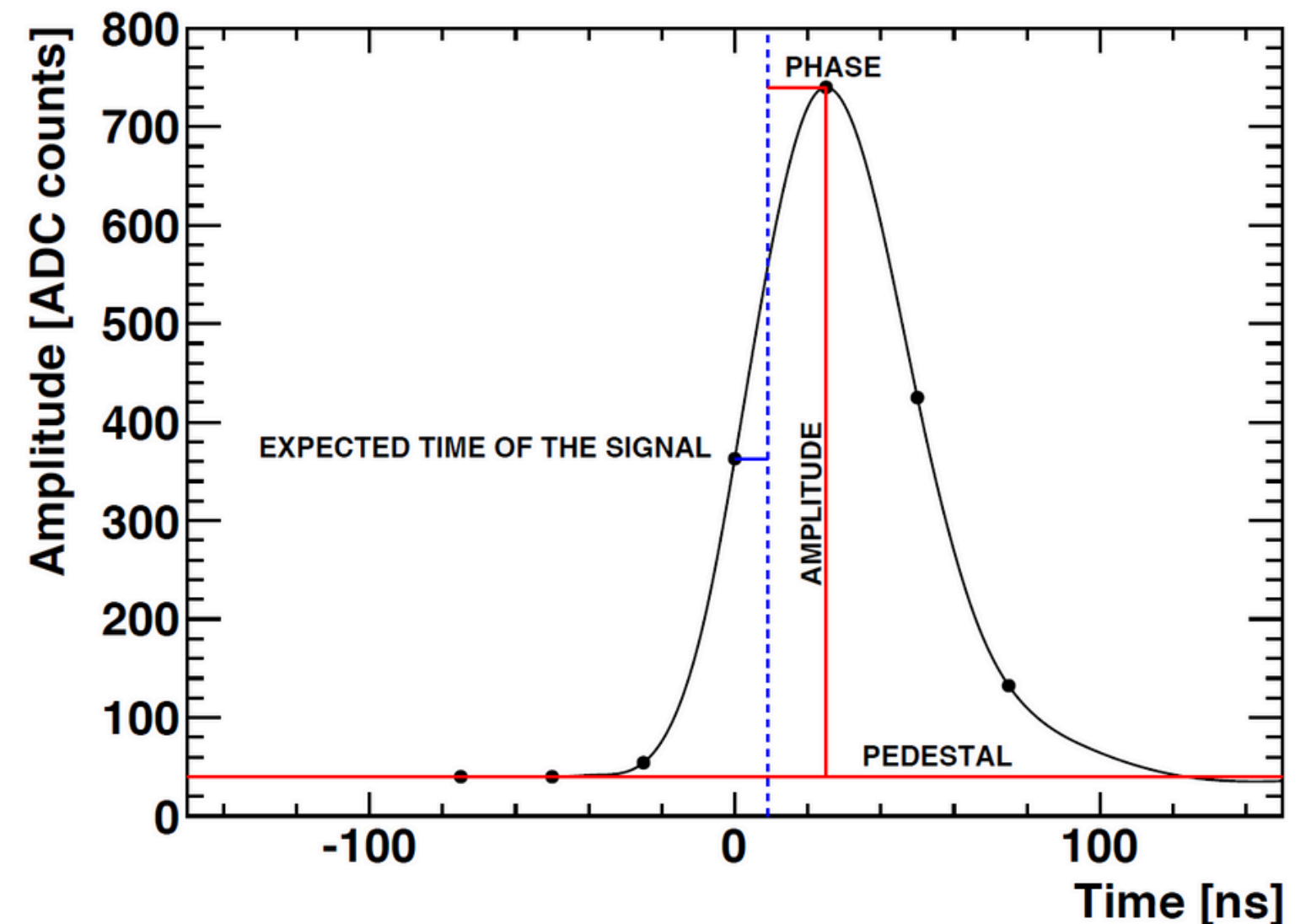  - **Fixed point** arithmetic is used

**Optimal Filtering Algorithm**

$$A = \sum_{i=1}^{n} a_i (S_i - p)$$

$$\tau = \frac{1}{A} \sum_{i=1}^{n} b_i (S_i - p)$$

$$\chi = \frac{1}{A} \sum_{i=1}^{n} |((S_i - p) - Ag_i)|$$

**DSP Online Algorithms for the ATLAS TileCal Read-Out Drivers**

Publisher: **IEEE** | Cite This | PDF

A. Valero ; J. Abdallah ; V. Castillo ; C. Cuenca ; A. Ferrer ; E. Fullana ; V. Gonzalez ; E. Higon ; J. Poveda ; A. Ruiz-Marti... **All Authors**
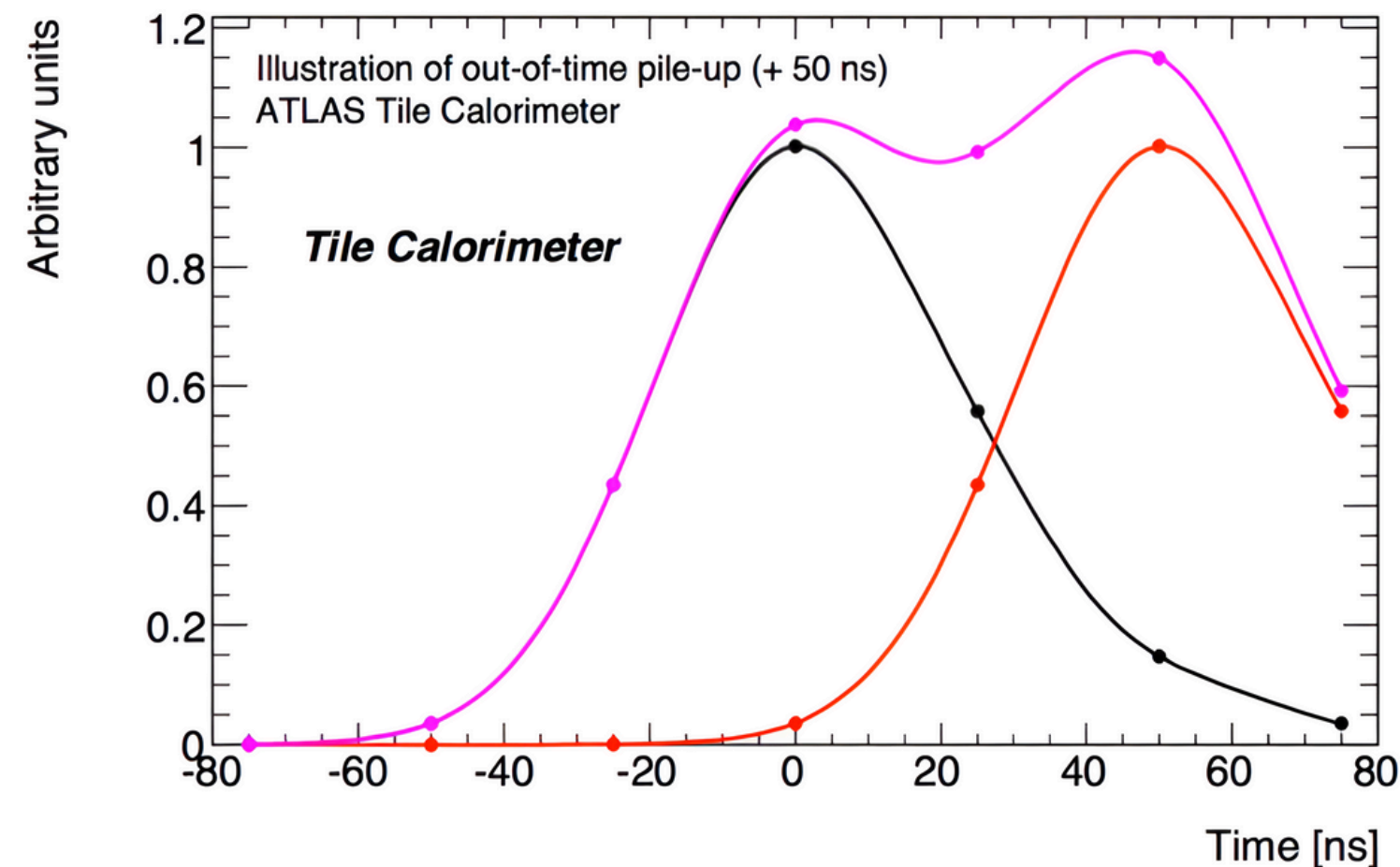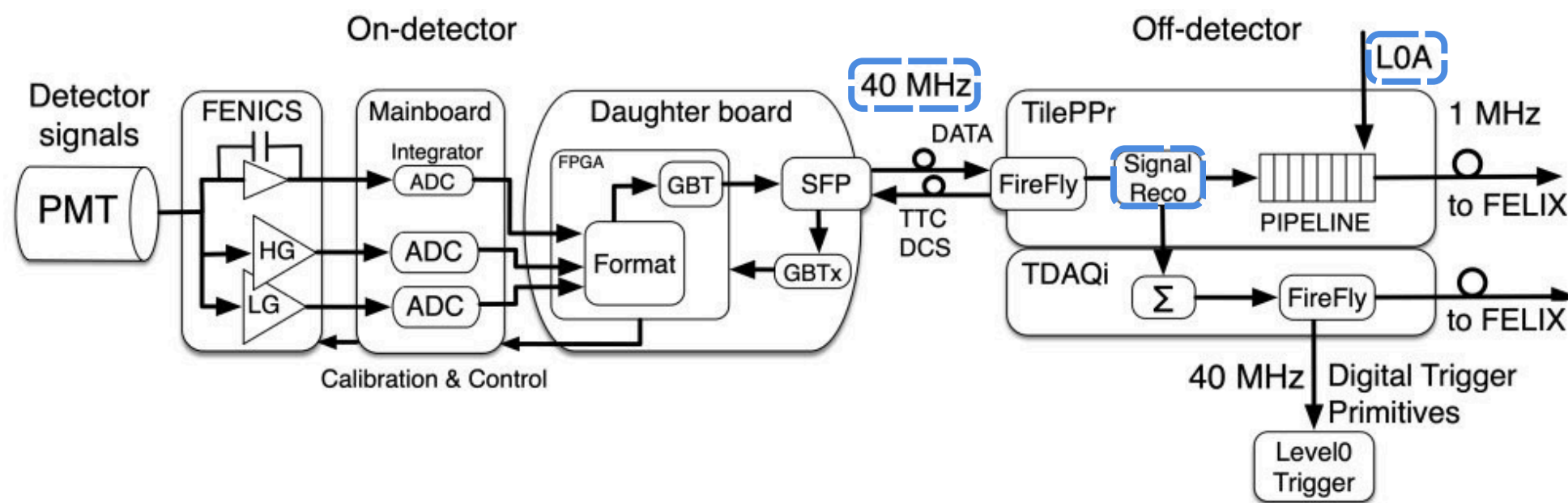
*Click on the image for reading more on the topic*

# 1 Introduction

## *HL-LHC calorimeter signal reconstruction*

- In the **HL-LHC**, signals will be **reconstructed for every BC** at **40 MHz** (25 ns) **before the trigger**
  - Signals need to be processed by **FPGAs** due to their **low and deterministic latency** for signal synchronization
  - Multiple **simultaneous signals** will produce **pile-up**
- There is a need for more sophisticated algorithms for signal reconstruction
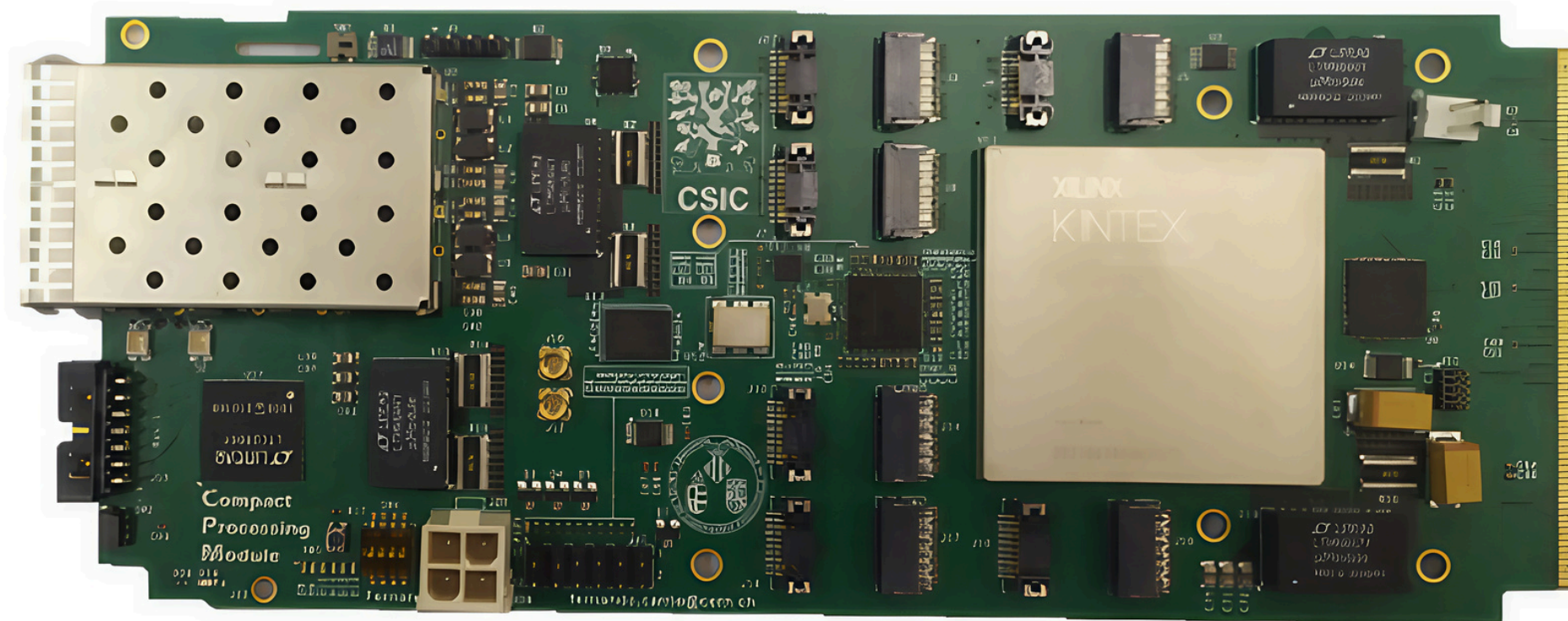  - **Deep learning** algorithms (**Neural Networks**)

# 1 Introduction

## *Real time processing*

- **Real time applications** fit better with **FPGAs**
- Algorithm replication must fit **area/occupancy**
- **Cycle accurate in latency** to interconnect modules

- **Real time** requirements:
  - Maximum algorithms in parallel: **162**
  - Maximum sample latency: **200 ns**
  - Sample frequency: **40 MHz**

**KU115 Chip resources:**

| | |
|---|---|
| FFs | 1326720 |
| LUTs | 663360 |
| Block RAMs (Mb) | 75.9 |
| DSP Slices | 5520 |



The PreProcessor module for the ATLAS Tile calorimeter at the HL-LHC

A. Valero*, F. Carrió, L. Fiorini, A. Cervelló, D. Hernandez and A. Ruiz Martinez

Instituto de Física Corpuscular, University of Valencia—CSIC, Valencia, Spain

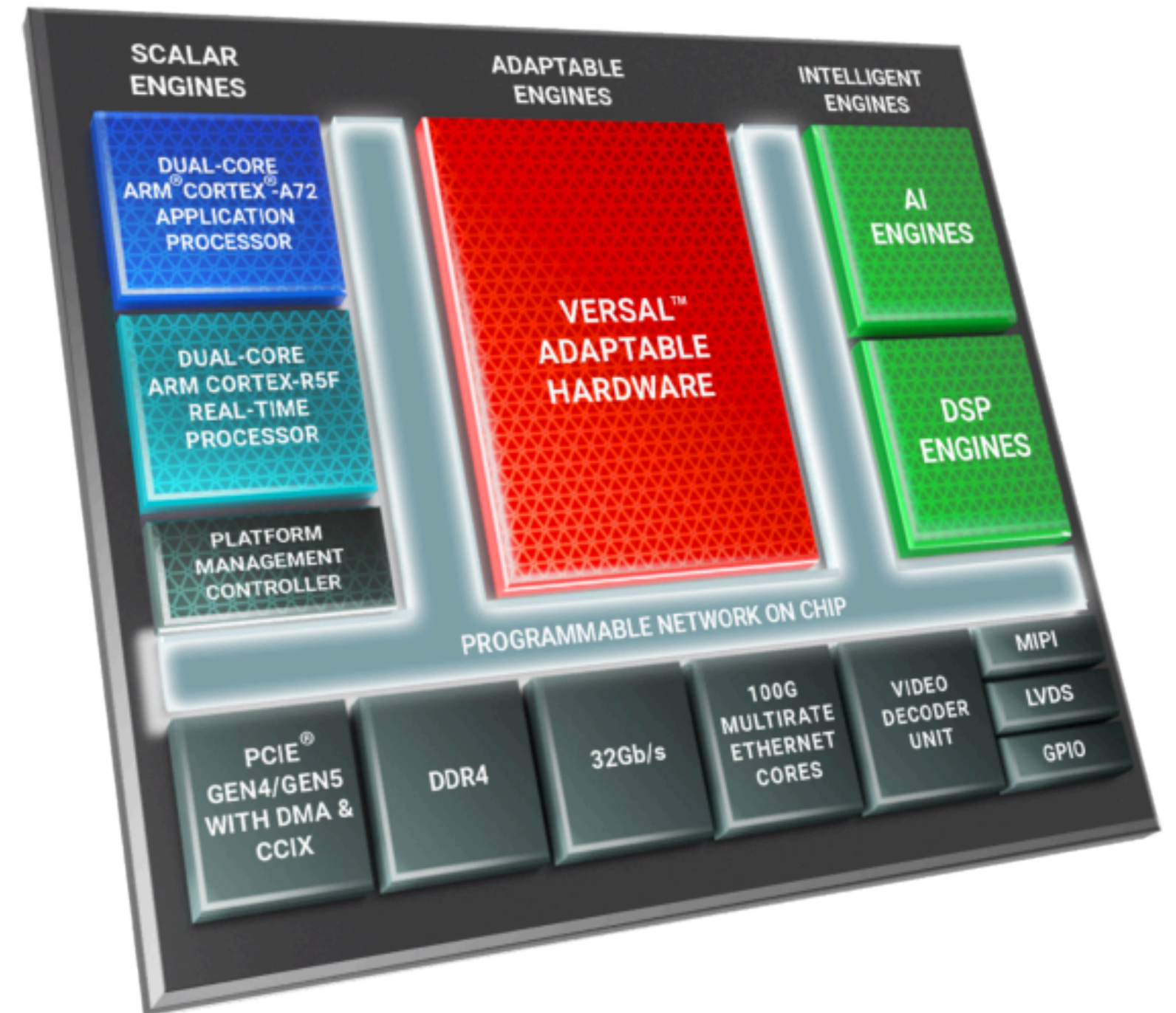*Click on the image for reading more on the topic*

# 2 Methodology

## *Versal ACAP device for algorithms test*

- **Versal ACAP** System on Chip (SoC) for **algorithm testing**
  - **PL - FPGA:** Algorithm hardware acceleration and data moving between memories
  - **PS - CPU:** Managing and control of the accelerators
  - **Memory:** Data buffering
  - **Connectivity:** External data reception and data transmission

**VC1902 Chip resources:**

| CPU | ARM A72 Dual core |
|---|---|
| | ARM R5F Dual core |
| Memory | OCM 256 KB |
| Connectivity | Ethernet x2 |
| | USB 2.0 x2 |
| | UART x2 |
| | SPI x2 |
| | I2C x2 |
| | CAN-FD x2 |

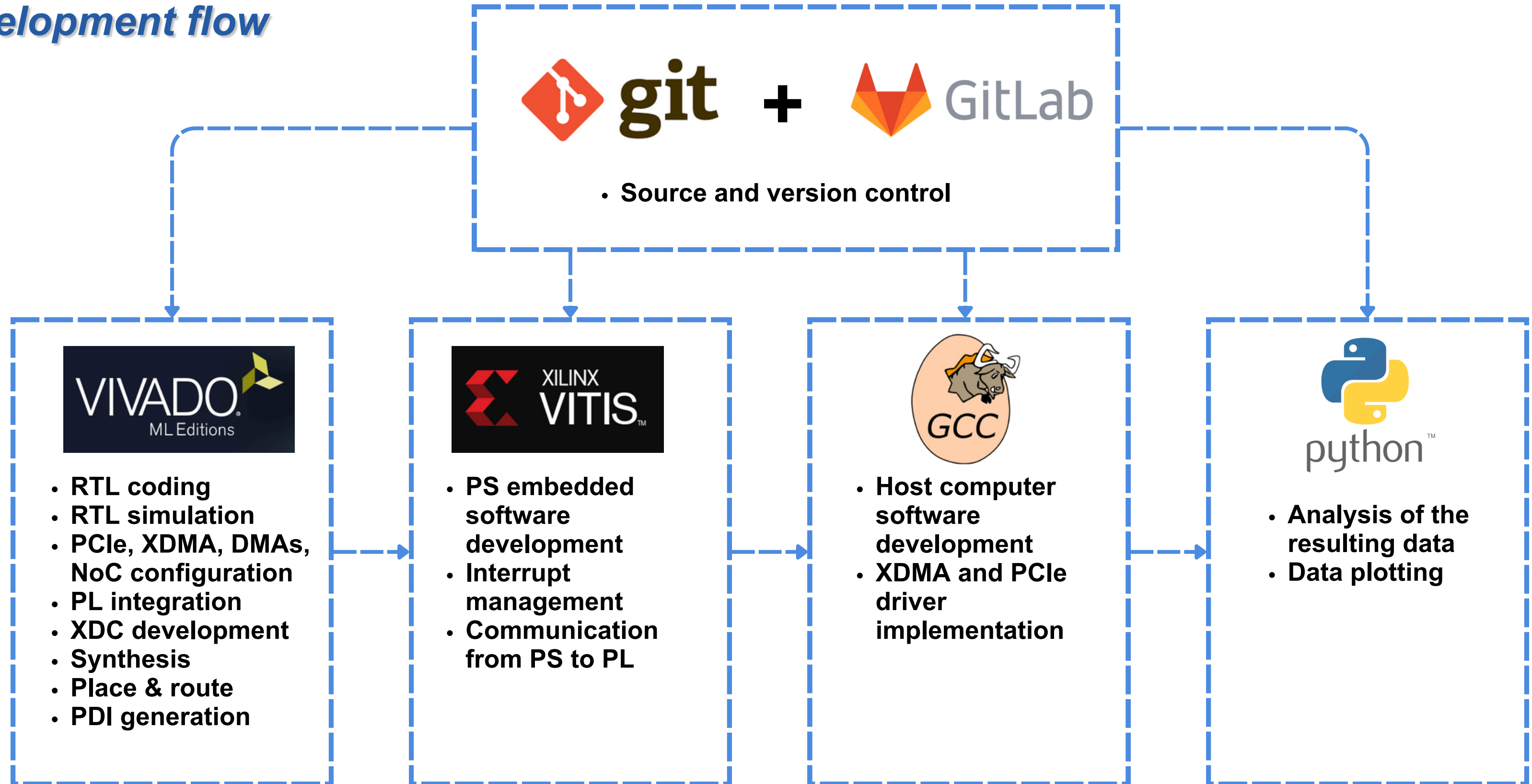| | |
|---|---|
| AI Engines | 400 |
| DSP Engines | 1968 |
| LUTs | 899840 |
| NoC Ports | 28 |
| DDR MC | 4 |
| PCIe – DMA | 1 |

# 2 Methodology

## *Setup development*

- Driver implementation in host CPU for communication with **XDMA**
- Driver implementation in device CPU for managing **internal DMAs**
- **NoC configuration** for internal communication
- **Interrupt system** development
- **Multiple cores** executing algorithms
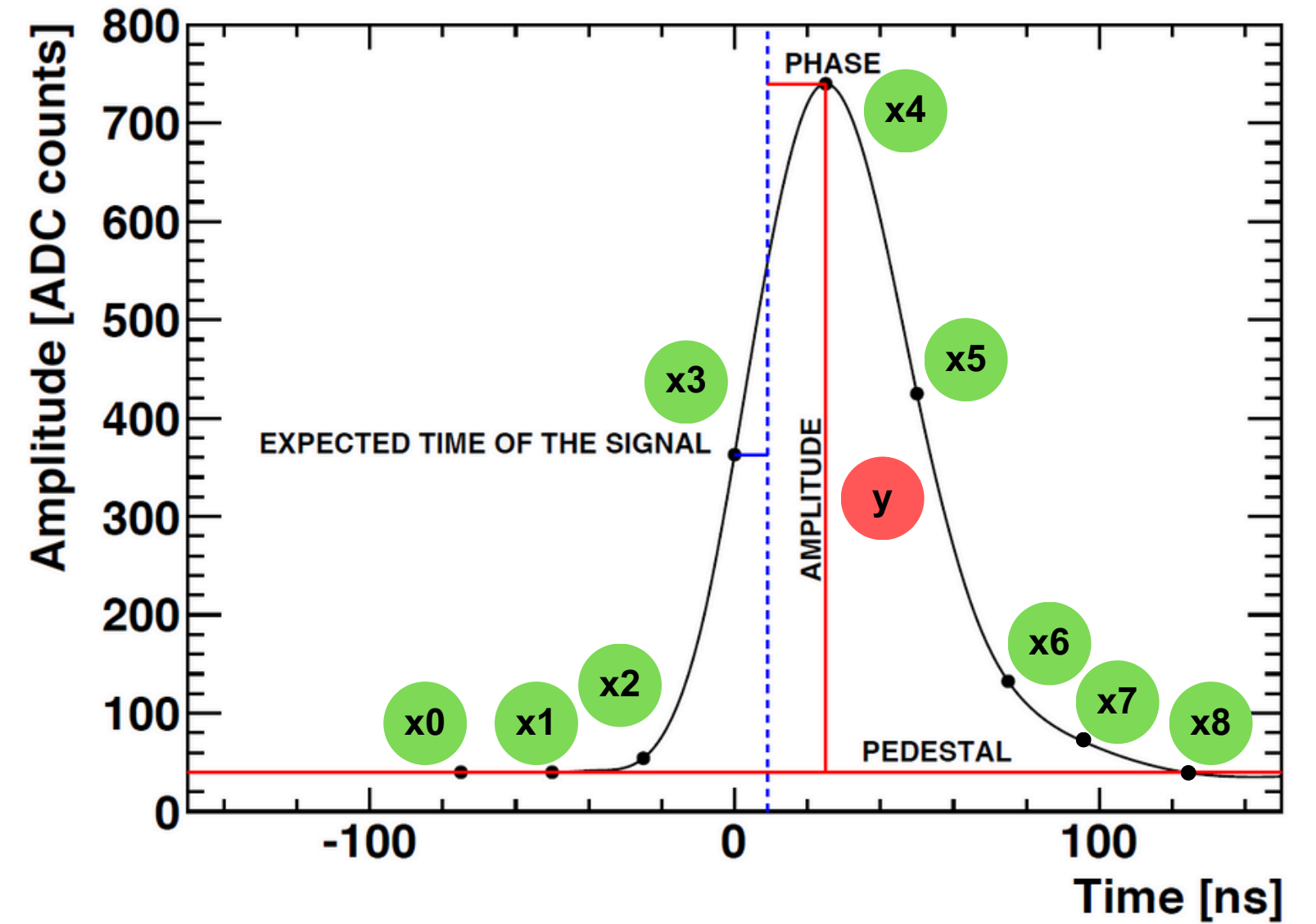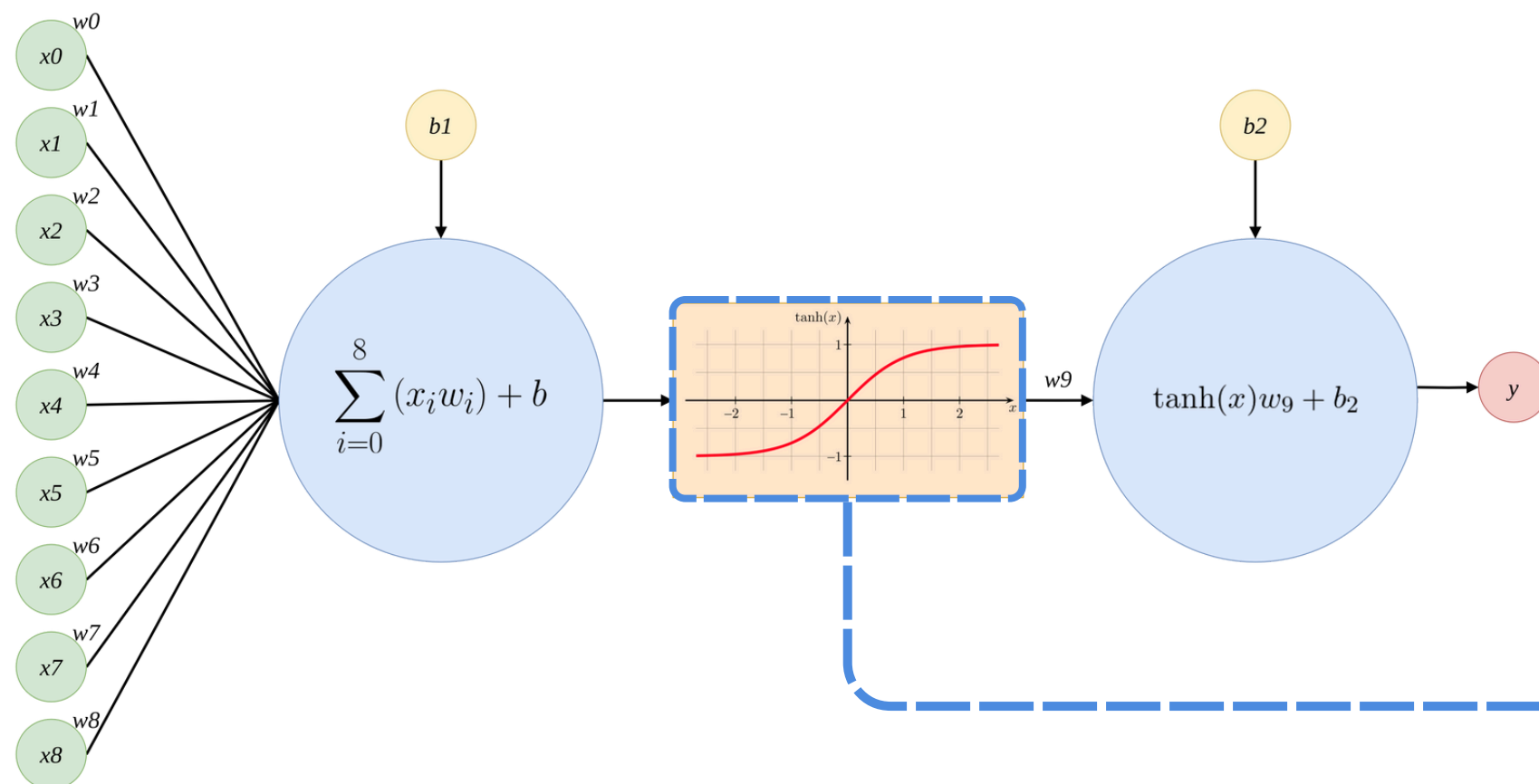


**PCIe Gen4 x8**

# 2 Methodology
## *Development flow*



- Source and version control

**VIVADO ML Editions**
- RTL coding
- RTL simulation
- PCIe, XDMA, DMAs, NoC configuration
- PL integration
- XDC development
- Synthesis
- Place & route
- PDI generation

**XILINX VITIS**
- PS embedded software development
- Interrupt management
- Communication from PS to PL

**GCC**
- Host computer software development
- XDMA and PCIe driver implementation

**python**
- Analysis of the resulting data
- Data plotting

# 2 Methodology

## Modified perceptron description

- Modified version of a **perceptron** with **2 neurons**
- **Hyperbolic tangent** as the activation function
  - Quantized and implemented in a Look Up Table
- Input: **9 BC** sliding window
- Output: **Amplitude** of the **central BC** in each window



FPGA implementation of a deep learning algorithm for real-time signal reconstruction in particle detectors under high pile-up conditions
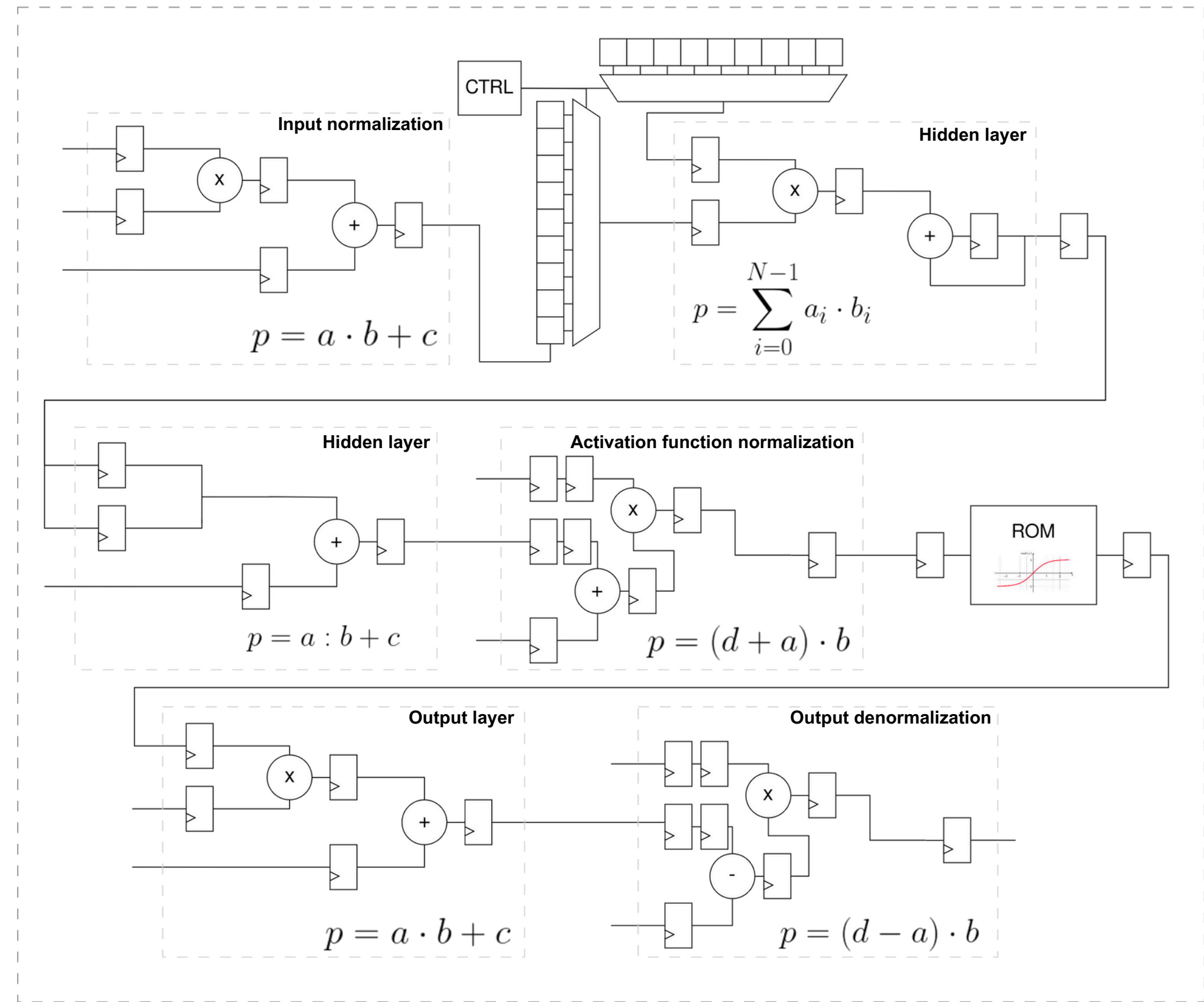
J.L. Ortiz Arciniega[1], F. Carrió[2] and A. Valero[2]

*Click on the image for reading more on the topic*

# 2 Methodology

## Modified perceptron RTL

- Written in **VHDL**
- Synthesized and implemented in **Vivado**
- **Fixed point** arithmetic
- Activation function **tanh(x) quantized** over **5000 values**
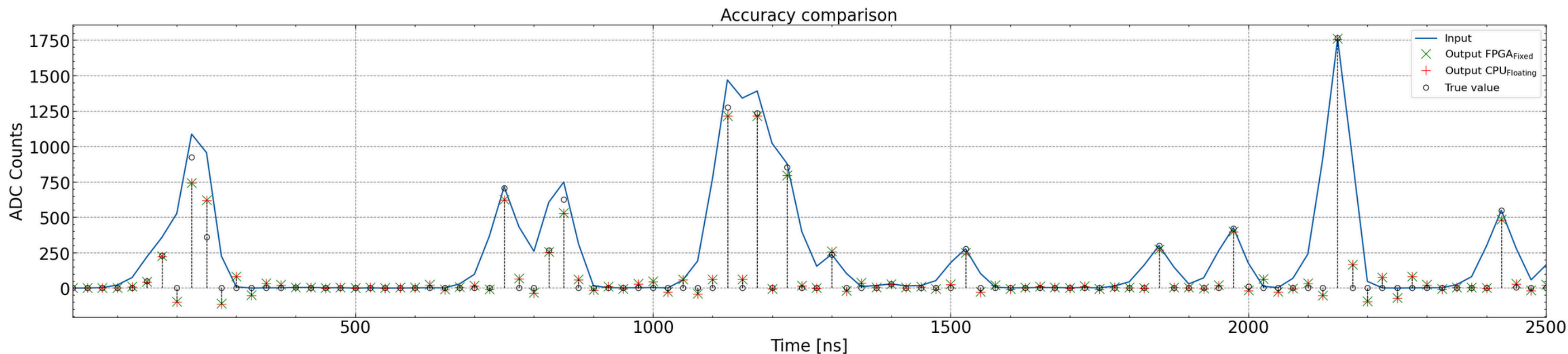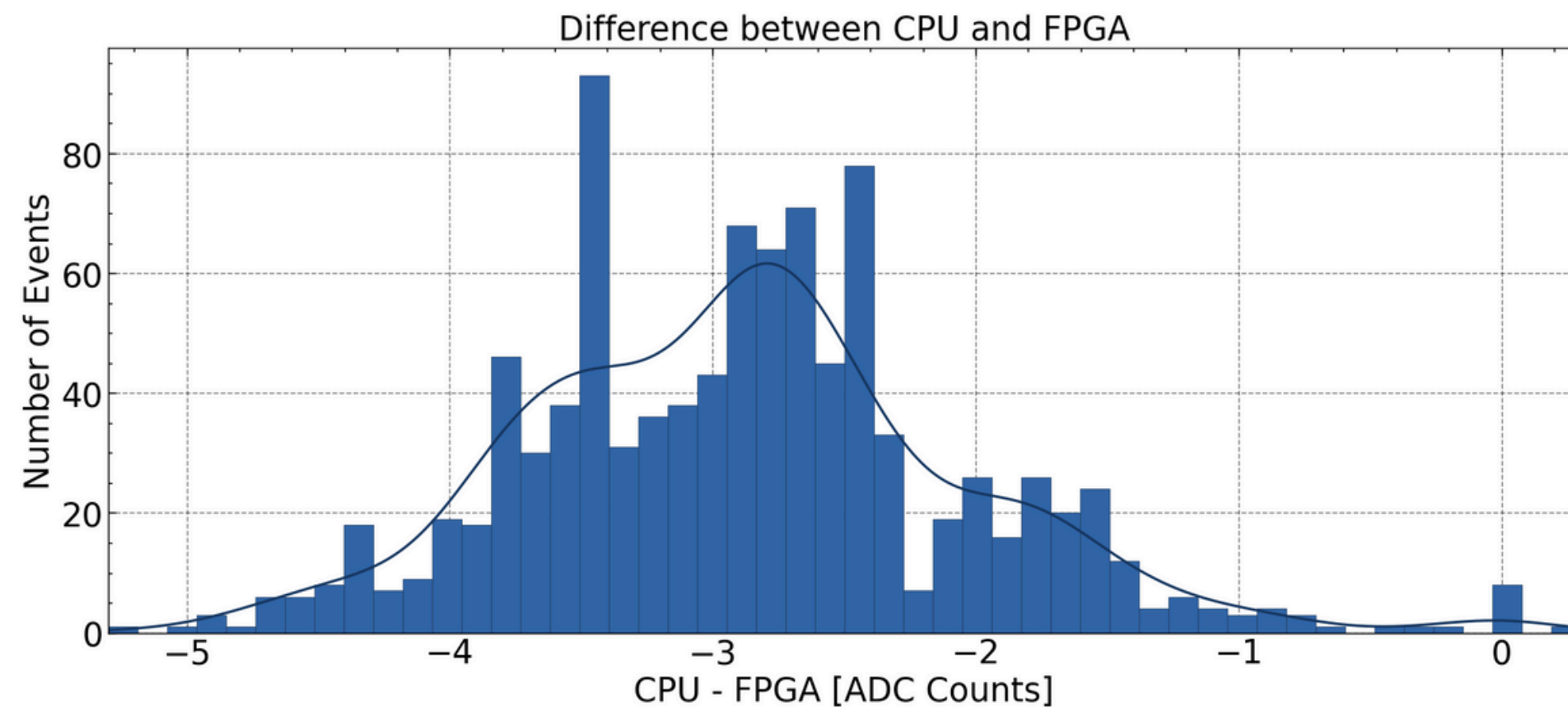- **Latency** of **34 clock cycles**

| VC1902 | | | |
|---|---|---|---|
| **Resource** | **Utilization** | **Available** | **Utilization (%)** |
| LUT | 72 | 899840 | 0.01% |
| FF | 203 | 1799680 | 0.01% |
| BRAM | 4 | 967 | 0.41% |
| DSP58 | 6 | 1968 | 0.30% |
| BUFG | 1 | 980 | 0.10% |

| KU115 | | | |
|---|---|---|---|
| **Resource** | **Utilization** | **Available** | **Utilization (%)** |
| LUT | 65 | 663360 | 0.01% |
| FF | 203 | 1326720 | 0.02% |
| BRAM | 3.5 | 2160 | 0.16% |
| DSP48E2 | 6 | 5520 | 0.11% |
| BUFG | 1 | 1248 | 0.08% |

CTRL

**Input normalization**

$$p = a \cdot b + c$$

**Hidden layer**

$$p = \sum_{i=0}^{N-1} a_i \cdot b_i$$

**Hidden layer**

$$p = a : b + c$$

**Activation function normalization**

ROM

$$p = (d + a) \cdot b$$

**Output layer**

$$p = a \cdot b + c$$

**Output denormalization**

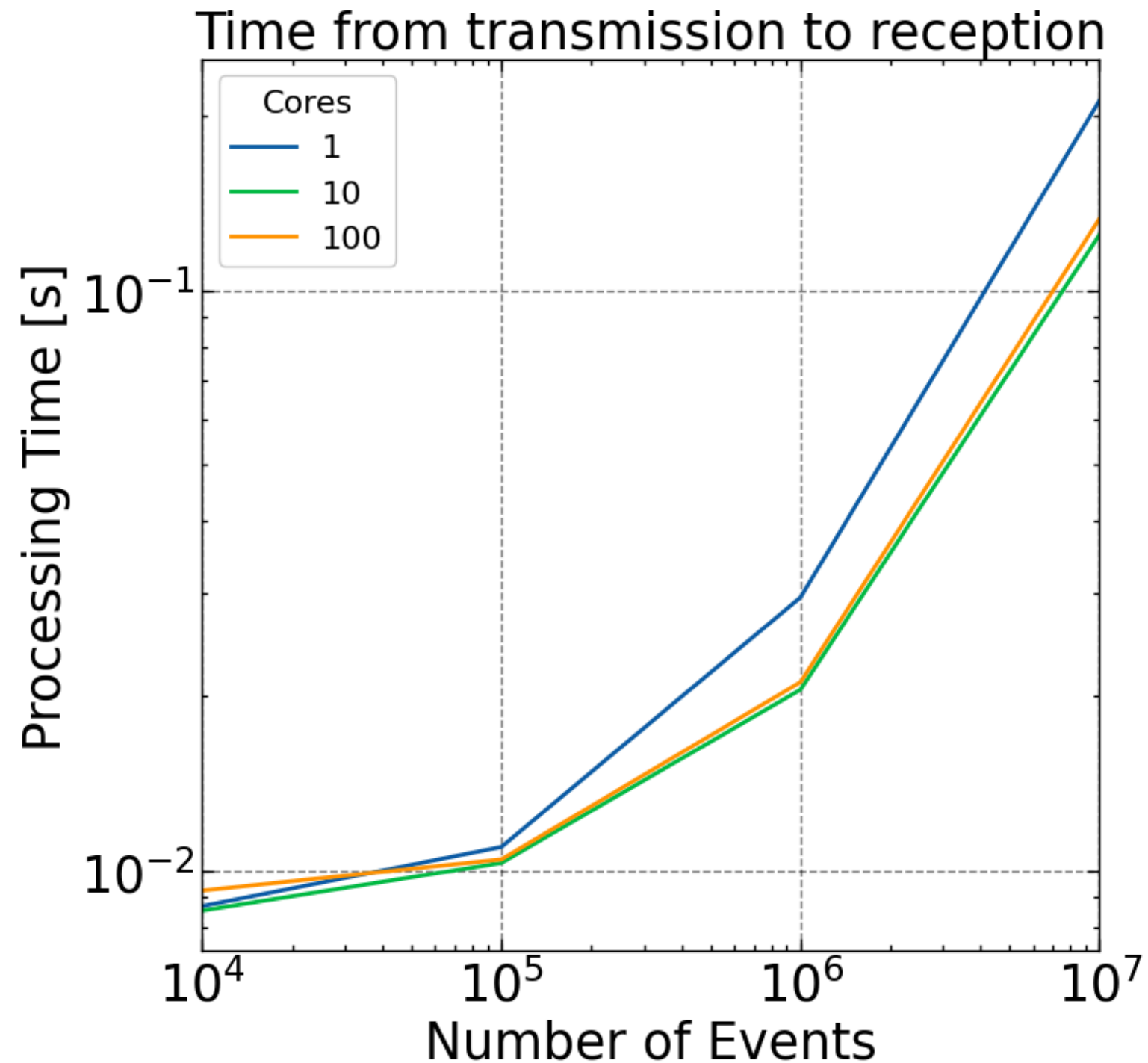$$p = (d - a) \cdot b$$

# 3 Results

## *Accuracy comparison*

- CPU (**Floating point**) vs. FPGA (**Fixed point**)
- Maximum difference: **5 ADC Counts**
- FPGA amplitude > CPU amplitud due to the **fixed point implementation**

# 3 Results

## *Multi-core implementation*



Time from transmission to reception
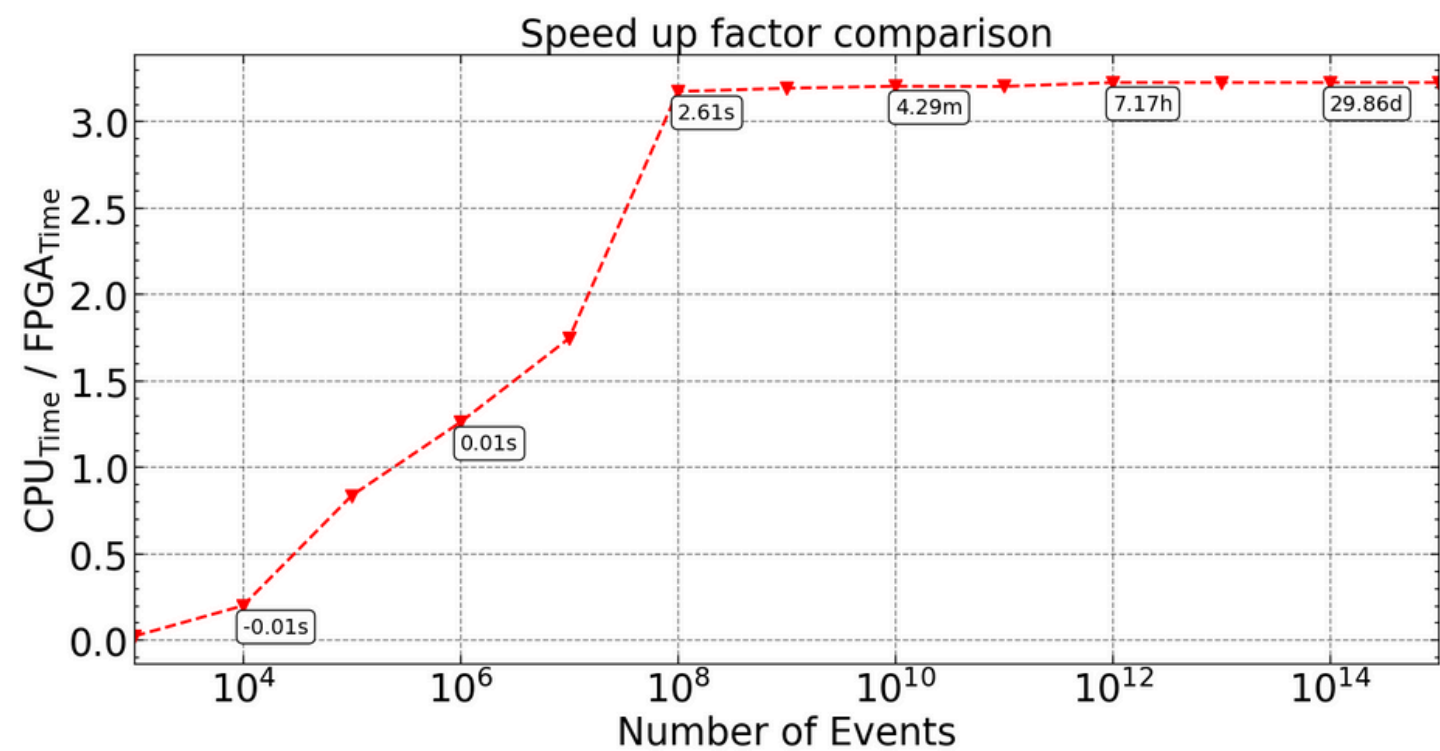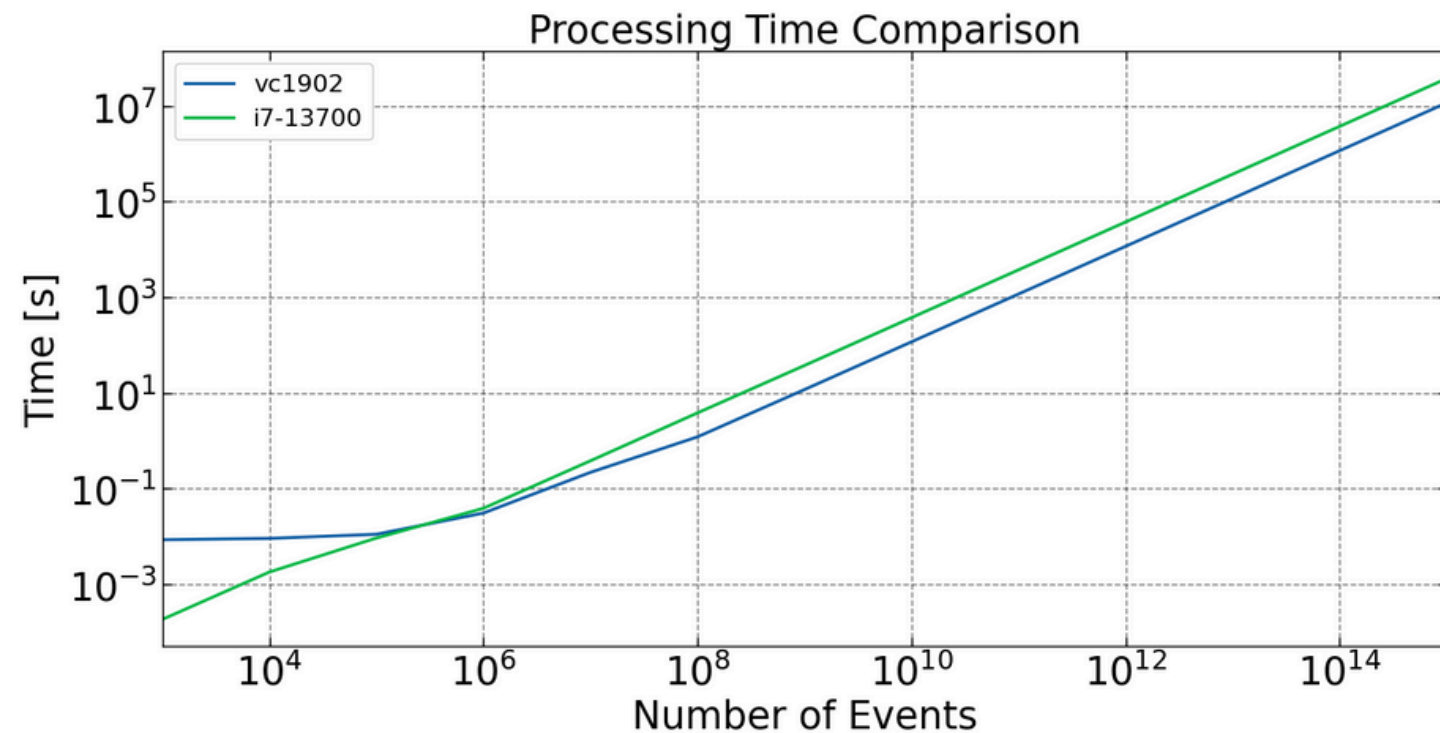
The **number of cores** is dependent of:
- NoC and DDR **bandwidth**
- NoC **OT** transactions
- PL **resources** used for each core (LUTs, FFs, BRAMs, DSPs, ...)

If **processing bandwidth** is greater than **NoC and DDR bandwidth**, **backpressure** is going to happen
- Source is **ready to send** data, **but** the consumer is **not read to receive**, so **data** is going to be **lost**

# 3 Results
## *Timing comparison*



- For **less than $10^6$ events**, the **CPU** is **better than** the **FPGA** due to the fixed minimum time for transmission and setup
- For **more than $10^6$ events**, the **FPGA** has a **better** performance **than** the **CPU**

- The **speed up factor** remains stable (**x3.2**) for **more than $10^8$ events**
- For **$10^{12}$ events**, the **FPGA is 7.17 hours faster than the CPU**

# 4 Conclusions
## *Summary and future work*

**Summary:**
- FPGA implementation of deep learning algorithms improves efficiency over traditional CPU
- Algorithm optimization for real time application is a trade-off between latency/throughput and area usage

**Future work:**
- Evaluate more complex deep learning algorithms for real time implementation
- Power consumption will be measured and monitored
- AI Engines utilization and optimization
- Optimization in terms of latency and power consumption

# Funding

# Versal ACAP Processing for HL-LHC Calorimeters Signal Reconstruction

**27th Conference on Computing in High Energy and Nuclear Physics (CHEP)**
**24th October, 2024**

*<u>Francisco Hervas</u>, Alberto Valero, Luca Fiorini, Hector Gutierrez*

# AI Engines

# AI Engines



- 400 AI Engine Tiles
- Frequency
  - 1.25 GHz working
  - 312.5 MHz transport
- Latency
  - Input net: 12 cycles
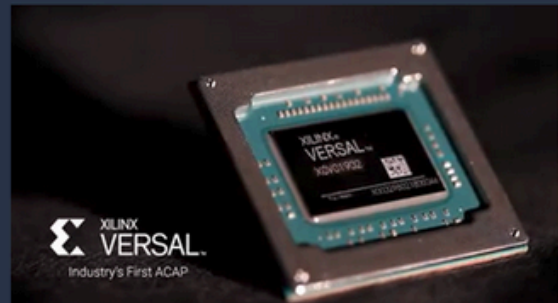  - Output net: 8 cycles

# BEAM System Controller

# Fixed point vs. Floating point



Floating-Point Format



Fixed-Point Format

# RTL Design